

CPExpertTM

DASD Component

Computer Management Sciences, Inc
6076-D Franconia Road
Alexandria, Virginia 22310-1756
voice: (703) 922-7027
fax: (703) 922-7305
www.cpexpert.com

CPExpert is a trademark of Computer Management Sciences, Inc.

This manual applies to the DASD Component of **CPEXpert™**, a proprietary software product of Computer Management Sciences, Inc., Alexandria, Virginia, United States of America.

The information in this document is subject to change. Comments or suggestions are welcome, and to the extent practicable, will be incorporated in revisions to this document. Please send all comments, suggested new rules, suggested changes to existing rules, suggestions for improvement to the software or documentation, or any other advice to:

Computer Management Sciences, Inc.
6076-D Franconia Road
Alexandria, Virginia 22310
(703) 922-7027 FAX: (703) 922-7305
www.cpexpert.com

DISCLAIMER

The advice, recommendations, or otherwise contained in this document represent information generally available in the public domain, as contained in vendor manuals, published in articles or papers, presented at professional conferences, or otherwise commonly accepted in the professional community. Neither Computer Management Sciences, Inc. nor its representatives make representations or warranties with respect to the applicability or application of any advice, recommendations, or otherwise, contained in this document or in any results from applying the CPEXpert software, to any particular computer system or computer installation.

TRADEMARKS

CPEXpert is a trademark of Computer Management Sciences, Inc. IBM, MVS/370, MVS/SP, MVS/XA, MVS/ESA, Enterprise System/3090, Netview, PR/SM, Processor Resource/System Manager, Hiperspace, and ES/9000 are trademarks of the IBM Corporation. SAS, SAS/OR, and SAS/STAT are trademarks of the SAS Institute Inc. MXG is a trademark of Merrill Consultants. MICS is a trademark of Legent Corporation.

COPYRIGHT INFORMATION

©Copyright 1991, Computer Management Sciences, Inc.
All rights reserved. Printed in the United States of America.

This licensed work is confidential and propriety, and is the property of Computer Management Sciences, Inc. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, or otherwise, without the specific written authorization of Computer Management Sciences, Inc.

Preface

Over the past several decades, computer technology advances have dramatically increased the processing power of central processors. Unfortunately, the improvements in I/O processing have not kept pace with the central processor advances. However, the data storage and accessing requirements of most installations continue to significantly increase as new applications are developed. The "DASD farm" often represents the largest investment in hardware in a large installation, and on-line applications require rapid access to data in order to meet response goals.

A variety of techniques have been implemented in an attempt to minimize the amount of I/O processing or to speed the access to data (dataspaces, virtual storage, expanded storage, and cached DASD are examples of techniques implemented to reduce the I/O requirement or to speed access). These techniques have had limited success in reducing the "I/O bottleneck" for most installations. Consequently, the DASD configuration still presents the largest potential for overall performance improvement and for performance improvement of specific applications. In fact, IBM has stated that *"over 75% of the problems reported to the IBM Washington System Center can be traced to some kind of I/O contention. Channel loading, control unit or device contention, data set placement, paging configurations, and shared DASD are the major culprits..."*¹

We believe that a typical Data Base Administrator can solve the more serious DASD performance problems with such information as:

- Which DASD devices offer the most potential for improved performance.
- Which DASD devices cause the most adverse effect on the performance of specific application systems (such as on-line applications).
- Which applications use these DASD devices and how much the applications use the DASD devices.
- Which applications interfere with on-line workload (such as TSO or CICS), and what is the effect of the interference.

The DASD Component of CPExpert is designed to identify DASD problems such as those listed above, and to provide the information necessary to solve the problems. Data Base Administrators can then focus on correcting the DASD problems identified.

¹MVS Performance Notebook, IBM publication GC28-0886

How to use this manual

This document describes how to use the DASD Component of CPEXpert to analyze major constraints to improved performance of your computer system. The manual is organized into six sections and one appendix.

Section 1 provides an introduction to the DASD Component. This section is organized into four chapters. Most of this section can be reviewed for general information.

- Chapter 1 provides a brief background of DASD performance problems and performance analysis.
- Chapter 2 provides an overview of the DASD Component of CPEXpert.
- Chapter 3 describes the sources of data used by the DASD Component to analyze DASD performance.
- Chapter 4 describes the performance data bases which CPEXpert can use to analyze performance.
- Chapter 5 briefly describes the types of analysis that the DASD Component performs.

Section 2 provides information on installing the DASD Component (please follow the instructions contained in the *CPEXpert Installation Guide* if CPEXpert has not been installed). The instructions in this chapter should be followed closely when installing the DASD Component, and may be required when you make changes to your I/O configuration. This section is organized into three chapters:

- Chapter 1 provides detailed instructions on how to install the modification to MXG necessary to collect information relating DASD use to specific applications (jobs, job steps, and service classes or performance groups). If you use MXG to create your performance data base, the instructions in this chapter will optionally be followed when you initially install the DASD Component.
- Chapter 2 provides detailed instructions on how to install the modification to NeuMICS necessary to collect information relating DASD use to specific applications (jobs, job steps, and service classes or performance groups). If you use NeuMICS to create your performance data base, the instructions in this chapter will optionally be followed when you initially install the DASD Component.
- Chapter 3 provides detailed instructions on how to (1) define workload categories, (2) associate service classes or performance groups with the workload categories, and (3) direct CPEXpert to perform analysis based upon the defined workload

categories. The workload definition is optional; you do not have to define any workloads. However, you must define workloads if you wish CPEXpert to analyze the DASD performance effects of contending workloads.

Section 3 describes how to provide guidance variables to the DASD Component using the CPEXPERT.SOURCE(DASGUIDE) PDS member. The instructions in this section will be important any time the guidance variables need to be changed. This section is organized into two chapters.

- Chapter 1 describes how to specify data selection and presentation guidance variables. The instructions in this chapter will be important if you wish to select specific measurement periods for analysis.
- Chapter 2 describes how to specify analysis control variables to guide the DASD Component in its analysis of system performance. The instructions in this chapter will be important if you wish to alter the defaults provided with the DASD Component.
- Chapter 3 describes how to exclude specific volumes from analysis.
- Chapter 4 describes how to select specific volumes for analysis.
- Chapter 5 describes how to analyze how well critical data sets meet response objectives.
- Chapter 6 describes how to analyze potential performance problems with VSAM data sets.
- Chapter 7 describes how to exclude specific VSAM data sets from analysis.

Section 4 provides information on executing the DASD Component. The instructions in this section should be followed closely when executing the DASD Component. The instructions in this section will be important each time the DASD Component is executed.

Section 5 describes DASD performance considerations. The information in this section should be appreciated before attempting to use the DASD Component. The section is organized into two chapters.

- Chapter 1 provides an overview of DASD performance considerations. This chapter describes the major components involved in a typical DASD I/O operation. The information in this chapter allows you to appreciate the importance of the areas analyzed by the DASD Component.

- Chapter 2 describes limitations and considerations involved in using SMF/RMF data to analyze DASD performance. The information in this chapter is essential in determining whether the analysis performed by the DASD Component is appropriate for your environment.

Section 6 describes how to use the DASD Component to analyze performance. The instructions in this section should be followed each time you execute the DASD Component.

Appendix A contains a detailed description of each rule that results in a finding based upon the DASD Component analyzing performance of your DASD. You may wish to briefly review the rules in this appendix to appreciate the problems that are encountered in different installations. However, it is not necessary to read all of the rules. It is necessary only to read the rules that are identified by the reports produced from the DASCPE Module.

Acknowledgments

Computer Management Sciences would like to acknowledge the following individuals. They played a significant role in the concept, design, development, testing, or documentation of CPExpert²:

John Ebner, Systemhouse

Alan Greenberg, Social Security Administration

Stan Meacham, MCI Corporation

Barry Merrill, Merrill Consultants

Philip Mugglestone, European Software Product Services, NV

Bryant Osborn, US Department of Transportation

John Peterson, MCI Corporation

Bernie Pierce, IBM Corporation

Fred Voth, JOSTENS Corporation

CPExpert would not exist as a product if it had not been for the personal inspiration, professional advice, technical knowledge, and encouragement from these individuals!

Additionally, Computer Management Sciences would like to acknowledge the following individuals. They played a significant role in the concept, design, development, and testing of the DASD Component of CPExpert:

V. Lee Conyers, US Department of Transportation

Carson Ho, UNISYS

Bryant Osborn, US Department of Transportation

²The affiliation of the individuals is shown as of the time CPExpert was developed. Some of the individuals are no longer associated with the organizations shown.

	<u>Page</u>
Changes	xi
 Section 1: Introduction	
Chapter 1: Background	1-1
Chapter 2: The DASD Component of CPEXpert	1-3
Chapter 3: Data Sources	1-5
Chapter 4: Performance Data Bases	1-6
Chapter 5: Types of Analysis	1-7
Chapter 5.1: Basic Analysis	1-7
Chapter 5.2: Expanded Analysis (Specific applications)	1-9
Chapter 5.3: Expanded Analysis (Specific data sets)	1-11
Chapter 5.4: Analysis of shared DASD conflicts	1-13
Chapter 5.5: Analysis of VSAM data sets	1-16
 Section 2: Installing the DASD Component	
Chapter 1: Installing the modification for MXG	2-3
Step 1: Install the DASDMXG code	2-3
Step 2: Modify MXG modules	2-3
Step 3: Modify MXG IMAC30DD module	2-6
Step 4(alternate): Modify MXG EXPDBINC module and EXPDBVAR module ..	2-8
Step 4(alternate): Modify the SAS job stream used to execute MXG	2-10
Step 5: Add CPEDASD DD statement to the JCL	2-10
Chapter 2: Installing the modification for NeuMICS	2-12
Step 1: Install the DASDMIC code	2-12
Step 2: Modify the sharedprefix.MICS.USER.SOURCE(#SMFEXIT)	2-12
Step 3: Modify prefix.MICS.SOURCE(DYSMFFMT)	2-14
Step 4: Add the CPEDASD DD statement to the JCL	2-17
Chapter 3: Defining workload categories	2-19
Chapter 4: Defining critical data sets	2-22
Chapter 5: Defining Multiple PDBs	2-25

	<u>Page</u>
Chapter 4: Defining critical data sets	2-22
Chapter 5: Defining Multiple PDBs	2-25
Chapter 1: Data Selection and Presentation Variables	3-2
Chapter 1.1: CONFIG variable	3-3
Chapter 1.2: CONFIGX variable	3-3
Chapter 1.3: DASDATES and DASTIMES variables	3-3
Chapter 1.4: DASDATEE and DASTIMEE variables	3-4
Chapter 1.5: DASDAT2S and DASTIM2S variables	3-4
Chapter 1.6: DASDAT2E and DASTIM2E variables	3-4
Chapter 1.7: MAXRULES variable	3-5
Chapter 1.8: SHIFT variable	3-5
Chapter 1.9: SYSTEM variable	3-6
Chapter 1.10: SYSTEMn variable(s)	3-7
Chapter 1.11: SYSPLEX variable	3-7
Chapter 1.12: SAS Output Delivery System	3-8
Chapter 1.13: VERBOSE variable	3-8
Chapter 2: Analysis Control Variables	3-10
Chapter 2.1: Number of devices to analysis: ANALYZE variable	3-11
Chapter 2.2: Exclude devices with low activity: DASDEXCP variable	3-12
Chapter 2.3: Analyze using response objectives: DASDSN variable	3-13
Chapter 2.4: Produce only I/O configuration: EVALDASD variable	3-13
Chapter 2.5: Excluding volumes from analysis: EXCLUDE variable	3-14
Chapter 2.6: Specifying data sets to list: LIST42DS variable	3-14
Chapter 2.7: Perform “loved one” analysis: LOVED1 variable	3-14
Chapter 2.8: Analyze all devices referenced by “loved one” applications: LOVEDALL variable	3-16
Chapter 2.9: Exclude reporting low-activity data sets: MIN42PCT variable ...	3-17
Chapter 2.10 List data for all RMF intervals: LISTALL variable	3-17
Chapter 2.11: Minimum I/O rate to analyze: MINIORT variable	3-18
Chapter 2.12: Minimum I/O rate to analyze: MINIOWT variable	3-18
Chapter 2.13: Minimum I/O response to analyze: MINRESP variable	3-18
Chapter 2.14: Number of volumes to report: REPORT variable	3-19
Chapter 2.15: Selecting volumes to analyze: SELECT variable	3-19
Chapter 2.16: Analyze shared DASD Conflicts - SHARED variable	3-19
Chapter 2.17: SMF Type 30 modification installed - TYPE30DD variable	3-20
Chapter 2.18: SMF Type 42 (Data Set Statistics)- TYPE42DS variable	3-21

	Page
Chapter 3: Excluding volumes from analysis	3-22
Chapter 3.1: EXCLUDE variable	3-22
Chapter 3.2: Defining volumes to exclude	3-22
Chapter 4: Selecting specific volumes for analysis	3-24
Chapter 4.1: SELECT variable	3-24
Chapter 4.2: Defining volumes to analyze	3-24
Chapter 5: Analyzing response objectives for critical data sets	3-26
Chapter 5.1: Analysis based on TYPE42DS	3-26
Chapter 5.2: Analysis based on TYPE14/15 and CPExpert modification	3-26
Chapter 6: Analyzing VSAM data sets	3-28
Chapter 6.1: Controlling analysis of VSAM: ANALVSAM variable	3-29
Chapter 6.2: Excessive Control Area splits: CASPLITS variable	3-30
Chapter 6.3: Percent direct to VSAM index component: DIRINDEX variable	3-31
Chapter 6.4: Excessive EXTENTS were allocated: EXTENTS variable	3-32
Chapter 6.5: Specifying LSR sequential domination: LSRSEQ variable	3-32
Chapter 6.6: Specifying maximum extents: MXEXTENT variable	3-33
Chapter 6.7: Specifying NSR direct access domination: NSRDIR variable	3-34
Chapter 6.8: Minimum VSAM open time: OPENTIME variable	3-35
Chapter 6.9: Specifying percent direct for CI size: PCTDIR variable	3-36
Chapter 6.10: Specifying percent sequential for CI size: PCTSEQ variable	3-37
Chapter 6.11: Excluding VSAM data sets: VSAMEXCL variable	3-37
Chapter 6.12: Specifying significant VSAM I/O activity: VSAMIO variable	3-38
Chapter 6.13: Summarizing VSAM activity: VSAMSMRY variable	3-39
Chapter 7: Excluding VSAM data sets from analysis	3-40
Chapter 7.1: VSAMEXCL variable	3-40
Chapter 7.2: Defining VSAM data sets to exclude	3-40

Section 4: Executing the DASD Component

Chapter 1: Executing the DASCPE Module	4-1
Step 1: Use TSO ISPF to change the "prefix" in the data set names	4-1
Step 2: Make any appropriate changes to the DASGUIDE Module	4-3
Step 3: Execute the DASCPE Module	4-3

	<u>Page</u>
Chapter 2: Executing the DAS1415 Module	4-4
Step 1. Use TSO ISPF to change the DD statements	4-4
Step 2. Execute the DAS1415 Module	4-5
Checklist for Executing the DASD Component, Mainframe	4-6
Checklist for Executing DASD Component, Personal Computer	4-7
Checklist for Performing Expanded Analysis	4-8

Section 5: DASD Analysis Factors

Chapter 1: Overview of DASD Performance Considerations	5-1
Chapter 1.1: IOSQ time	5-1
Chapter 1.2: PEND time	5-2
Chapter 1.3: DISC time	5-4
Chapter 1.4: CONN time	5-6
Chapter 1.5: OTHER time	5-7
Chapter 2: RMF Data Analysis Considerations	5-9
Chapter 2.1: SMF information	5-9
Chapter 2.2: Data Averages	5-11

Section 6: Using the DASD Component

Chapter 1: Prepare guidance for the DASD Component	6-1
Chapter 2: Actions on a daily basis	6-2
Step 1: Execute the DASCPE Module	6-2
Step 2: Review the output from the DASCPE Module	6-2
Chapter 3: Actions on a weekly or monthly basis	6-3

Appendix A: Description of Rules

Page**Exhibits**

2-1	MXG IMACINTV Module, with CPEXpert Modification	2-14
2-2	MXG EXTY30U4 Module, with CPEXpert Modification	2-15
2-3	MXG IMAC30DD Module, with CPEXpert Modification	2-17
2-4	Sample MXG SAS Job Stream	2-21
2-5	Normal MICS.USER.SOURCE(#SMFEXIT) before modification	2-23
2-6	Normal MICS.USER.SOURCE(#SMFEXIT) after modification	2-23
2-7	Sample MICS.SOURCE(DYSMFFMT), with CPEXpert Modification	2-27
2-8	Sample CPEDASD DD statement in MICS JCL	2-27
2-9	Sample display of CPEXPERT.USOURCE(DASGUIDE) Module	2-29
2-10	Sample display of CPEXPERT.USOURCE(DASGUIDE) Module	2-32
2-11	Sample Display of CPEXPERT.USOURCE(DASGUIDE) Module	2-35
3-1	Data Selection and Presentation Variables	3-2
3-2	Analysis Guidance Variables	3-11
3-3	Excluding specific volumes from analysis	3-23
3-4	Selecting specific volumes for analysis	3-25
3-5	VSAM analysis guidance variables	3-28
3-6	Excluding VSAM data sets from analysis	3-40
4-1	Job Control Language to execute the DASCPE Module	4-2
4-2	Job Control Language to execute the DAS1415 Module	4-4
5-1	Major Components of DASD I/O Operations	5-1

Changes

CPEXpert Release 13.2:

The main changes to the DASD Component for CPEXpert Release 13.2 are to:

- C Update the discussion of DASD performance problems and specific findings of the DASD Component, to include discussion of FICON infrastructure implications, Parallel Access Volume (PAV), and cached device issues.
- C Specifically identified all findings that relate to “legacy” systems (e.g., 3380 devices attached to 3990-2 controllers) so readers will not be confused about discussions that do not apply to more modern environments.
- C Provide better analysis of cache controller features, operation, and performance implications.
- C Provide better analysis of device DISC time (including DISC time caused by physical channel activity).
- C Provide better analysis of device PEND delay time (including PEND delay time caused by cache controller activity and miss hits).
- C Create an approach whereby the configuration model created by the DASD Component can be retained in a specific library (rather than the SAS WORK library), and provide an option to process from this library rather than re-create the configuration model. Retaining the configuration model allows capacity planners to access the configuration model for capacity planning purposes. Providing the option to process the retained configuration model allows performance analysts to examine problem areas repeatedly without incurring the overhead of creating the configuration model.
- C Add physical channel type (e.g., ESCON, FICON Bridge, FICON Native, etc.) and physical channel activity to the configuration model so future analysis can detect performance problems with channels based on type of channel.
- C Add RMF Cache Controller statistics to the configuration model (only for those devices that are selected for detailed analysis). These cache controller statistics allow current analysis of cache controller performance problems (and facilitate expended analysis in future releases of CPEXpert).
- C Provide the following new rules:
 - C DAS131 PEND time was caused by channel busy |
 - C DAS132 PEND delay time was caused by director port busy |
 - C DAS133 PEND delay time was caused by controller busy delays |

- C DAS134 PEND delay time was caused by device busy delays
- C DAS135 PEND delay time was caused by other delays
- C DAS160 Disconnect was major cause of response delay

CPExpert Release 13.1:

The main changes to the DASD Component for CPExpert Release 13.1 are to:

- C A new feature allows users to specify target response times for specific data sets. The DASD Component will use information from TYPE42DS records to detect data sets that miss a specified response objective. The devices on which the data sets reside will then be analyzed to identify performance problems.
- C Variables have been added so a user can specify that CPExpert should ignore devices that are below a minimum I/O response, below a minimum I/O rate, or below a minimum total I/O wait time; and to suppress listing RMF intervals for devices that have no problems in the interval.
- C The CPExpert modification to MXG or neuMICS has been modified to revised the "MAXDASD variable value exceeded:" message. This message can be produced by the CPExpert code that extracts device information from SMF TYPE30 records (as either MXG or neuMICS processes the SMF data). The previous message indicated that the MAXDASD variable was too low and some devices were ignored by CPExpert. One implication of this message was that user tended to increase the MAXDASD value to very large numbers. Since this value controls array sizes in many modules of the DASD Component, a large value for MAXDASD caused unnecessary overhead. With the new approach, the code lists the jobs that have excess devices. A user can decide whether these jobs are sufficiently important to warrant the increased execution time of the DASD Component.
- C The logic that processes TYPE42DS has been revised to significantly increase execution efficiency of the DASD Component.
- C Add the ability to select up to 10 systems individually for analysis. Until Release 131, a user had the options of analyzing data for all systems in the performance data based, analyzing data for a specific sysplex (in case the performance data base contained data for more than one sysplex), or analyzing data for a specific system in the performance data base. With Release 13.1, up to 10 systems can be individually selected for analysis.

- C Enhance the options provided with the SAS Output Delivery System (ODS). With Release 13.1, users who exercise the SAS ODS feature for creating CPEXpert output can optionally create the output as a PDF file, which can be emailed to other users. Additionally, users can optionally specify a STYLE feature for either HTML or PDF output, if they have a preferred STYLE for HTML or PDF output. The optional links that are available with the HTML have been revised; SAS at some user sites did not create the HTML output in the “standard” way, and the CPEXpert code that inserted links into the HTML output did not work properly. I have revised the code to place the links into the output as the output is created, rather than attempting to place the links into the final HTML output created by SAS.
- C CPEXpert now specifies OPTIONS COMPRESS=N; to override any site specification for file compression. Experiments have shown that CPEXpert code runs significantly faster (using much less CPU time) if compression has been turned off.

CPEXpert Release 12.2:

The main changes to the DASD Component for CPEXpert Release 12.2 are to:

- C Enhance the DASD Component to analyze problems with VSAM data sets¹. This analysis is a partial automation of the analysis and guidance given in IBM’s *VSAM Demystified* Redbook, SG24-6105. The VSAM analysis is performed if a MXG performance data base exists, and if SMF Type 42 (Data Set Statistics) and SMF Type 64 files exist in the MXG performance data base. These are MXG files TYPE42DS and TYPE64, respectively. The following rules have been added to the DASD Component as a part of the VSAM analysis:

DAS600: Excessive Control Area (CA) splits occurred

DAS604: Excessive secondary extents were allocated

DAS605: Excessive extents were used and secondary allocation unit was small

DAS606: Primary or Secondary allocation unit was small

DAS607: VSAM data set is close to maximum number of extents

DAS610: Relatively small CI size was used for sequential processing

DAS611: Relatively large CI size was used for direct processing

DAS612: Relatively large CI size was used for mixed processing

¹Thanks to Glenn Bowman (Wakefern Food Corporation, NJ), Joan Kelley (IBM Poughkeepsie, NY), and John Cothran (IBM Dallas, TX) for providing VSAM test data.

DAS620: The number of data buffers should be increased

DAS621: The number of index buffers should be equal to index levels

DAS622: The number of index buffers should be more than STRNO value

DAS625: NSR was used, but a large percent of the access was direct

DAS635: LSR was used, but a large percent of the access was sequential

- C Options are provided to analyze VSAM data sets residing only on “poorly performing” devices, analyze all VSAM data sets, analyze only VSAM data sets (suppressing the normal “worst devices” analysis), or suppress analysis of VSAM data sets.
- C An option is provided to suppress analysis of DASD altogether, but simply create the model of the I/O configuration, and associated device/controller/channel activity. This option can be useful for reporting I/O activity for capacity planning.
- C A report optionally is produced when common analysis is not performed because data sources are not available (or CPExpert has not been advised that the data sources are available). For example, the report will alert you to missing application analysis if the modification to MXG or MICS has not been made so CPExpert has application data available.
- C A report optionally is produced if the DASD report is excessively large, and suggestions are given about how to reduce the size of the report.
- C The DASD Component is updated to support z/OS Version 1 Release 4.

CPExpert Release 12.1:

The main changes to the DASD Component for CPExpert Release 12.1 are to:

- C Completely revise the Component to eliminate the requirement that users provide IOCP macros from which CPExpert built a model of the I/O configuration. IBM now provides sufficient information in standard SMF records so the DASD Component can automatically create I/O configurations using SMF Type 70, Type 73, Type 74, Type 75, and Type 78CF. With Release 12.1 of CPExpert, users will no longer need to provide IOCP macros to the DASD Component, since CPExpert can obtain the information automatically.
- C Eliminate all documentation related to IOCP macros and other documentation for user-provided input related to the I/O configuration.

- C Provide data set access information related to those data sets (managed by DFSMS) residing on volumes with poor performance.
- C Include a detailed analysis of PEND time, when PEND time is a major cause of performance problems.
- C Include the ability to analyze more than one sysplex in a single execution of the DASD Component
- C Enhance the output to provide a sysplex view of DASD volumes with poor performance, regardless of whether a specific system is being analyzed.
- C Update the Component to support z/OS Version 1 Release 3.

CPEXpert Release 11.2:

The main changes to the DASD Component for CPEXpert Release 11.2 are to:

- C Update the Component to support z/OS Version 1 Release 2.
- C Add support for the SAS Output Delivery System (ODS) feature, to enable optional web access of CPEXpert reports (this new option was suggested by **Harald Seifert** of HUK, Coburg, Germany).
- C Add optional links in CPEXpert reports (if the SAS ODS feature is invoked), that link rule output to CPEXpert documentation for the rules produced.
- C Revise the entire DASD Component User Manual to correct administrative errors and to make to document easier to use.

CPEXpert Release 11.1:

The main changes to the DASD Component for CPEXpert Release 11.1 are to update the Component to support z/OS Version 1 Release 1.

CPEXpert Release 10.2:

The main changes to the DASD Component for CPEXpert Release 10.2 are to update the Component to support OS/390 Version 2 Release 10.

CPEXpert Release 10.1:

The main changes to the DASD Component for CPEXpert Release 10.1 are to:

- C Update the Component to support OS/390 Version 2 Release 9.
- C Revise the DASD Component User Manual to correct Section 2, Chapter 4 (Installing the modification for MICS).
- C Revise the DASD Component User Manual to remove the requirement that the GENPARMS module be executed.

CPEXpert Release 9.2:

The main changes to the DASD Component for CPEXpert Release 9.2 are to update the Component to support OS/390 Version 2 Release 8.

CPEXpert Release 9.1:

The main changes to the DASD Component for CPEXpert Release 9.1 are to:

- C Update the Component to support OS/390 Version 2 Release 7.
- C Correct documentation to reflect existing logic to handle more than 4 paths to a device.
- C Add a description of techniques that are used to generate an IOCP input data set, using the Hardware Configuration Definition (HCD) component of MVS.
- C Add documentation of Rule DAS102 and Rule DAS202 to describe the rules that have been a standard part of the output (but had been left out of the User Manual).
- C Revise the description of installing the modification for MXG if you use the MXG BUILDPDB or use SAS/ITSV to process SMF data.
- C Revise output based on suggestions from users.
- C Correct code based on errors reported by users.

CPEXpert Release 8.2:

The main changes to the DASD Component for CPEXpert Release 8.2 are to:

- C Update the Component to support OS/390 Version 2 Release 6.

- C Correct code based on errors reported by users.

CPExpert Release 8.1:

The main changes to the DASD Component for CPExpert Release 8.1 are:

- C Added code to support Workload Definitions for MVS (Goal Mode), when exercising the “Loved One” analysis option. The effect of this change is to allow users to specify Service Classes for Goal Mode processing, rather than specify Performance Groups as done in pre-Goal Mode operation.
- C Modified the DASDMXG module to include Report Classes for MVS (Goal Mode). Additionally, added code in all relevant modules to support Report Class processing.

Section 1: Introduction

This section provides a brief overview of DASD performance evaluation, provides an overview of the DASD Component of CPEXpert, describes the sources of data used by the DASD Component to analyze system performance, and describes the performance data bases that the DASD Component can use.

Chapter 1: Background

From a global view, DASD analysis can be viewed from two perspectives;

- How the DASD space is allocated (how much space is allocated versus how much space is free), how the allocated space is used (whether the allocated space is actually used by the files), and how frequently DASD files are referenced by applications (whether the files should be on-line or perhaps stored off-line). DASD performance analysis from this perspective typically relates to the cost of maintaining the on-line storage represented by the DASD devices.
- How the DASD configuration and files effect overall system performance or effect the performance of individual applications. DASD performance analysis from this perspective deals with the day-to-day operational performance implications of the DASD configuration and files.

There are numerous products that examine DASD performance from the first perspective mentioned above. CPEXpert does not depend upon these products, since a basic philosophy of CPEXpert is to use information normally available in a performance data base. Consequently, the DASD Component does not analyze DASD space allocation.

The DASD Component focuses on DASD performance from the second perspective: whether performance constraints exist in the basic DASD configuration and to what degree particular applications impact the performance of other, more important, applications.

Beginning with Release 12.2 of CPEXpert, the DASD Component also analyzes problems with VSAM data sets. This analysis is a partial automation of the analysis and guidance given in IBM's *VSAM Demystified* Redbook, SG24-6105.

Very sound DASD analysis advice was given by Friesenborg¹: "It seems that the main thing to do is to avoid the outrageous; after that the subject doesn't warrant much attention." The major objective of DASD analysis is to identify and solve the "outrageous" or serious problems. The DASD Component of CPEXpert attempts to detect the serious problems, identify the causes of the problems, and suggest solutions to the problems.

¹Friesenborg, S. E., IBM Publication GG22-9370

DASD performance constraints can be corrected by a variety of actions at different levels:

- At the file level (e.g., change file block size, move files to different location)
- At the device level (e.g. upgrade device)
- At the controller level (e.g., change the number of paths, upgrade controller,, add cache)
- At the system level (e.g., reconfigure paths)
- At the operations level (e.g., schedule conflicting applications to different times)
- At the applications level (e.g., change storage accessing requirements or patterns).

The DASD Component provides information and recommendations regarding which actions are suitable for specific problems in your environment.

Chapter 2: The DASD Component of CPExpert

The DASD Component was implemented with the basic philosophy that:

- Most Data Base Administrators are able to address a limited number of DASD problems each day. The DASD Component was designed to identify the most serious DASD problems (both from the perspective of the overall DASD configuration and from the perspective of critical applications).
- Data Base Administrators are most interested in solving recurring DASD problems. The DASD Component was designed to identify the devices, files, and applications responsible for or experiencing recurring DASD problems.

With this design philosophy, the DASD Component (1) analyzes DASD information available in standard SMF/RMF records to select the most serious problems, and (2) identifies the devices and applications responsible for the problems or applications experiencing performance degradation because of the problems.

The guidance provided to the DASD Component allows a wide variety of options in the analysis process. These options allow users to define workload categories, to exclude specific DASD volumes from analysis, and to specify thresholds to limit the analysis.

These features and the analysis are provided by numerous SAS modules. These SAS modules (1) shape DASD performance and utilization data for detailed analysis by other modules, (2) determine the workload and analysis categories, (3) evaluate the data to assess potential areas of unacceptable performance, and (4) report the results from the evaluation. These modules are loaded and controlled by the central DASD Component of CPExpert (titled DASCPE).

- **Shape Data.** The majority of "processing" in the DASD Component is accomplished by numerous modules whose function is to combine, sort, summarize, correlate, and prepare the data for analysis. These modules extract and process relevant data elements from SMF records. Additionally, the modules incorporate a variety of queuing models to perform preliminary analysis of the data.
- **Determine workload and analysis categories.** A significant optional feature of the DASD Component is its ability to allow the specification of workload categories (e.g., TSO, CICS, Batch, etc.) and to focus analysis with respect to these workload categories. Additional controls provide the flexibility to exclude DASD volumes from consideration.
- **Evaluate Data.** The evaluation of the data is accomplished by rules whose purpose is to determine whether DASD performance was acceptable, and to analyze the DASD devices with unacceptable performance. These rules determine the devices with major problems, identify the probable causes of the problems, and (if expanded analysis is performed) identify the applications most likely responsible for the problems.

- **Report Results.** The DASD Component reports the results from the evaluation, in a variety of reports at the device or application workload level.

Each significant finding is described in Appendix A of this document. The description summarizes the finding, lists predecessor findings, discusses the rationale for the finding, and recommends action.

- The summary presents a short description of the finding.
- The predecessor findings are listed so you can follow the line of reasoning leading to a particular rule being executed.
- The discussion describes as much as necessary of the operation of the computer system as it relates to the particular finding. The purpose of the discussion is to explain the reasoning behind the finding. If appropriate, the discussion might refer you to related discussions in the DASD Component User Manual, or in a User Manual of another CPExpert component (e.g., the User Manual for the WLM Component may be used as a reference to avoid repeating a detailed discussion).
- The recommendations suggest possible actions that should be considered based on the findings. In many cases, multiple possible actions are listed. You must determine which actions should be taken (this determination is based upon the suitability of the actions to your own environment, the financial implications of the action, and the "political" acceptability of the action.)

Chapter 3: Data Sources

CPEXpert analyzes the performance of the DASD configuration based upon data from three sources:

- **RMF information.** The basic analysis performed by the DASD Component relies on standard SMF Type 70(series) information. The DASD Component analyzes Type 70, Type 72, Type 74, Type 75, and Type 78CF information to build a model of the I/O configuration in the sysplex, and to analyze performance of devices that performed poorly relative to other devices in the configuration.
- **SMF information.** The DASD Component of CPEXpert optionally includes a modification to MXG or to NeuMICS (using standard exit facilities). The modification allows extraction of SMF Type 30(Data Definition) information related to the individual DASD devices used by job steps . This information allows the DASD Component to perform "expanded" analysis, analyzing DASD performance based upon critical workloads.

Additionally, the DASD Component optionally extracts data from SMF Type 42 (Data Set Statistics) to report information about the data sets residing on devices with poor performance.

- **Guidance information.** Guidance information is contained in CPEXPERT.USOURCE(DASGUIDE). This PDS member contains information that allows you to specify specific time periods in which CPEXpert should analyze data, which performance groups (for Compatibility Mode) or service classes (for Goal Mode) are associated with different workload categories, what level of detail should be pursued in the DASD analysis, which volumes should be excluded from analysis, device path and caching characteristics, and other variables to control the operation of the DASD Component.

Chapter 4: Performance Data Bases

The "raw" RMF data contained in the SMF Type 30(DD) and Type 70(series) records must be translated into SAS format and placed into a performance data base before CPEXpert can use the information. There are two ways in which these records can be placed into a performance data base:

- **MXG performance data base.** The performance data base can be created by Merrill's Expanded Guide (MXG) software. MXG is provided by Merrill Consultants, Dallas, Texas. MXG provides a low-cost mechanism by which installations can create and maintain a performance data base.
- **NeuMICS performance data base.** The performance data base can be created by the MVS Integrated Control System (NeuMICS). NeuMICS is provided by Computer Associates, Vienna, Virginia. The NeuMICS SMF (or Base) Component creates and maintains a comprehensive performance data base from SMF/RMF information.

The flexibility to use either of the above performance data bases is due to the fact that CPEXpert is implemented in the SAS language. SAS provides a powerful macro facility, both with respect to macro coding and with respect to macro variable names.

CPEXpert generally uses SAS macro variable names when referring to an element of information in the SMF records. CPEXpert uses the SAS "%LET" statements to define the macro variables as MXG variable names or NeuMICS variable. These "%LET" statements are contained in unique variable definition modules for MXG and NeuMICS.

Thus, the same CPEXpert software can be executed against either of the performance data bases, by only invoking the proper definition module. The SAS %LET statements in the definition module automatically cause CPEXpert to refer to the proper MXG or NeuMICS data elements.

Chapter 5: Types of Analysis

The DASD Component performs two types of analysis of DASD performance and constraints to improved performance: **basic** analysis and **expanded** analysis. Additionally, the expanded analysis can analyze DASD performance from the perspective of specific application workloads or from the perspective of specific data sets. With CPExpert Release 12.2, the DASD Component analyzes potential performance problems with VSAM data sets. This section describes these types of analysis.

Chapter 5.1: Basic Analysis

The philosophy of the basic analysis is to identify the devices that have the most potential for improvement, and to analyze these devices in detail. The method used to identify a candidate set of devices is quite simple:

- CPExpert computes the average device response time for each **type** of device in the configuration, for each RMF measurement interval. The logic computes the average device response by type of device, since better performance would be expected from cached devices (for example) than from non-cached devices. This method essentially assesses the performance of each device against the performance of similar devices in the configuration.
- Devices that exceed the average device response time for their device type in any RMF measurement interval are selected as candidates for improvement. The rationale is that improvement efforts should not be directed at devices that provide better than average response. Thus, the candidate set of devices to analyze consists of those that provided worse than average response.
- The I/O rate of each candidate device is weighted by its response time. The result is a measure of the relative performance improvement **potential** of each device that provided worse than average response, from an overall system view. For example, consider two devices in a device type having an average I/O response of 20 milliseconds:

Device A: I/O rate = 30 I/O operations per second
Device response = 25 milliseconds
Weighting factor = $30 * 25 = 750$

Device B: I/O rate = 5 I/O operations per second
Device response = 40 milliseconds
Weighting factor = $5 * 40 = 200$

In the above example, CPExpert would select Device A as having the most overall potential for improvement, even though its per-I/O device response was not as bad as the device response of Device B.

This logic totally ignores the potential situation where Device B holds critical data sets (an index, for example) and thus has more effect on overall system performance than the I/O rate indicates. Unfortunately, CPExpert is examining measurement data and has no way to assess relative importance of devices. (In the basic analysis, CPExpert **does** provide an option by which you can select specific devices, and CPExpert will analyze performance constraints of only those devices.)

- A model is built of the I/O configuration (using Type 70, Type 73, Type 74, Type 75, and Type 78CF information). Measurement data (channel, controller, device information) is processed against this model. CPExpert uses queuing models to compute the likely delays at various physical or logical parts of the configuration. The result from this analysis associates probable delays with the specific devices.
- The device(s) with the overall most potential for improvement is selected for detailed analysis as to likely causes.
- The results from the above analysis are reported via rules described in Appendix A.

Chapter 5.2: Expanded Analysis (Specific applications)

CPEXpert can perform an expanded analysis of DASD performance from the perspective of specific applications. This section describes how this analysis is accomplished.

- A modification is made to the MXG or NeuMICS software creating your performance data base. This modification normally uses standard user exits², adds an insignificant amount of processing time, and uses a small amount of DASD space.
- The modification normally is invoked as MXG or NeuMICS process the SMF Type 30(Interval) records.

For MXG, the modification **could** be invoked as MXG processes the SMF Type 30(Job Step Termination) records if you do not collect SMF Type 30(Interval) records. Using Job Step Termination records is discouraged since incomplete analysis may result (the job steps may span many RMF measurement intervals and it may be impossible to correlate the data recorded by RMF with lengthy job steps).

- The modification records the system, job, job step, performance group or service class, and summary information about each job step's use of DASD devices by device number. The information is recorded only at the device level (rather than at the DD name level) and only required information is retained. Consequently, the records are small and require a small amount of DASD space.
- You define specific applications (or "loved one" applications) to CPEXpert. Specific workload categories are defined (for example, TSO or CICS), and you define the performance groups or service classes assigned to the workload categories.
- CPEXpert processes the data acquired using the modification to MXG or NeuMICS described above, and matches application use of devices against SMF Type 74 device information. The resulting information is application-specific, in that it relates device use to the applications using the devices.
- CPEXpert then identifies the devices used by "loved one" workloads, computes the average device response from this set of devices, and selects all devices exceeding the average response as candidates for analysis.
- The normal basic analysis is performed of all devices resulting from the above selection criteria. The result is an analysis from the perspective of specific applications.

²A single one-line modification is required to NeuMICS code in the NeuMICS DYSMFFMT module. NeuMICS does not currently provide an exit that can be used to add a data set to the output from NeuMICS as it analyzes SMF data. Such an exit does exist for the NeuMICS CICS Component (the _USRSDKP exit), and a request has been made that a similar exit be provided for the NeuMICS BASE Component.

- Additionally, CPExpert can identify non-critical applications causing performance problems for the "loved one" workload. For example, CPExpert can identify all non-critical applications using a significant amount of a path if CPExpert discovers that path contention results in poor performance to "loved one" applications.

Chapter 5.3: Expanded Analysis (Specific data sets)

CPEXpert can perform an expanded analysis of DASD performance from the perspective of specific data sets. The analysis is done in one of two ways: (1) analyzing data set response based on information from TYPE42DS statistics or (2) analyzing data based on information in TYPE14/15 and using the CPEXpert modification to MXG or NeuMICS as the SMF Type30 records are processed.

- **Analysis based on TYPE42DS.** This method can provide comprehensive information, without requiring the modification to MXG. This method **is not applicable** to NeuMICS, since NeuMICS does not retain sufficient information related to SMF Type 42 (Data Set Statistics).

MXG creates TYPE42DS from the SMF Type 42 (Subtype 6 - data set I/O statistics) records created by SMS. The TYPE42DS file contains I/O access characteristics information, at the data set level. CPEXpert extracts data set information and data set response statistics from TYPE42DS, and compares these to the response objectives you have made in USOURCE(DASGUIDE) about critical data sets. When any data set response time exceeds the specified objective, the DASD Component selects the data set and the volume it resides on for detailed analysis.

- **Analysis based on TYPE14/15 and CPEXpert modification.** This method is a bit more involved, but can be used regardless of whether data sets are managed by SMS, and this method can be used regardless of whether you use MXG or NeuMICS to create your performance data base.

With this method, the DASD Component analyzes SMF Type 14/15 records to extract data set names that correspond to the critical data sets that you have identified in USOURCE(DASGUIDE). The CPEXpert modification to MXG or NeuMICS extracts DD information as SMF Type 30 records are processed. The DASD Component then correlates the data set information with the DD information, to determine whether critical data sets exceed the specified response objective. When any data set responses time exceeds the specified objective, the DASD Component selects the data set and the volume it resides on for detailed analysis.

With either of the above methods, CPEXpert matches the data set information against SMF Type 74 device information, and identifies each device that has an average I/O response time exceeding the response objective of any data set residing on the device. The resulting information is specific to devices that have I/O response times exceeding a response objective.

The I/O rate of each device identified above is weighted by its response time. The result is a measure of the relative performance improvement potential of each device that provided I/O response exceeding the response objective for some data set on the device.

The normal basic analysis is performed of all devices resulting from the above selection criteria. The result is an analysis from the perspective of specific data sets.

Additionally, CPExpert creates a report describing the I/O activity of all data sets that resided on devices having an I/O response worse than the response objective for the data set.

Chapter 5.4: Analysis of shared DASD conflicts

CPEXpert can perform an analysis of conflicts between DASD shared between systems or MVS images. The analysis performed by CPEXpert is not intended to identify an isolated performance problem. Rather, CPEXpert attempts to identify those problems that **continually** cause shared DASD performance problems. Shared DASD analysis is invoked by specifying **%LET SHARED = Y;** in **USOURCE(DASGUIDE)**. Shared DASD analysis is an option, because more processing is required to perform shared DASD analysis.

CPEXpert performs the following processing if you have indicated that an analysis of potential conflicts between shared DASD should be performed:

- CPEXpert determines whether the "worst" devices selected for detailed analysis are shared with another system. If so, CPEXpert performs an analysis of potential conflicts caused by shared DASD.
- CPEXpert identifies other systems that reference the "worst" device. This identification is accomplished by analyzing the SMF Type 74 data in the performance data base relating to all other systems. The SMF Type 74 data contain the VOLSER for each device referenced. CPEXpert simply selects SMF Type 74 information for the systems that reference the VOLSER of the "worst" device. This information is retained for more detailed analysis about potential conflicts.

There is a potential (but very unlikely) problem with this method of identifying devices shared between systems. Multiple systems in the performance data base could use the same VOLSER to identify different devices. This could happen if the devices were not shared between systems.

For example, suppose that CPEXpert had identified PAGE01 as the "worst" device. Several system in the performance data base could reference VOLSER PAGE01, but the devices with VOLSER PAGE01 could be unique to each system. CPEXpert would assume that all references by other systems to PAGE01 applied to the "worst" device being analyzed. The references could apply to a totally different device, and the other systems might not even share DASD with the system being analyzed.

If this should be a problem (that is, if the DASD Component reports shared DASD conflicts with systems that do not share the device being analyzed), simply ignore the analysis produced by CPEXpert³.

³We do not feel that this problem will be common. It is described only to alert you to a potential incorrect analysis. If any user encounters this problem and it becomes annoying, code can be implemented to allow users to identify specific systems that share DASD with the system being analyzed. At present, this option seems to add unnecessary complexity to the user options.

- Once CPEXpert has identified all systems that reference the "worst" device, CPEXpert analyzes the DASD I/O characteristics of these systems with respect to the "worst" device. CPEXpert will produce Rule DAS300 to list statistics relating to potential conflicts, by system, by volume, and by RMF measurement interval.

The shared DASD analysis is accomplished when you perform the following steps:

- You collect and process SMF data for each system to be analyzed. The SMF data must be placed into your standard performance data base. This step is not specific for CPEXpert; it is the normal processing you would do. The step is listed first just for completeness.
- **You should ensure that the system clocks on all systems are set at roughly the same time.** CPEXpert will examine the SMF data contained in your performance data base, based upon the SMF time associated with the measurement data. Shared DASD problems between systems will be detected based on this SMF time. If the SMF times are significantly different, the analysis might not properly identify conflicts or conflicting applications.

You should not worry that the system clocks be set at **exactly** the same time. The analysis performed by CPEXpert is not intended to identify an isolated performance problem. Rather, CPEXpert attempts to identify those problems that **continually** cause shared DASD performance problems. If a shared DASD problem is continual, the problem will be reflected in multiple SMF recording intervals. Consequently, it is not essential that the SMF times be exactly matched between systems.

- You specify **%LET SHARED = YES;** in USOURCE(DASGUIDE) to tell CPEXpert to perform analysis of shared DASD.

CPEXpert performs the following processing if you have indicated that an analysis of potential conflicts between shared DASD should be performed:

- CPEXpert determines whether the "worst" device selected for detailed analysis is attached to a control unit shared with another system. If so, CPEXpert performs an analysis of potential conflicts caused by shared DASD.
- CPEXpert identifies other systems that reference the "worst" device by analyzing the SMF Type 74 data in the performance data base relating to all other systems. The SMF Type 74 data contain the VOLSER for each device referenced. CPEXpert simply selects SMF Type 74 information for the systems that reference the VOLSER of the "worst" device. This information is retained for more detailed analysis about potential conflicts.
- Once CPEXpert has identified all systems that reference the "worst" device, CPEXpert analyzes the DASD I/O characteristics of these systems with respect to the "worst" device.

- If "expanded" analysis is being performed, CPExpert can identify the applications on System B that reference the "worst" device.

SMF Type 30 records do not contain the VOLSER in the DD section. However, CPExpert correlates the information acquired from other records to relate the VOLSER to the device number. While extracting data from the SMF Type 74 information in the performance data base, CPExpert retains the device number associated with the devices. This device number is matched against the device numbers in the SMF Type 30 records to identify the VOLSER associated with the device.

CPExpert thus can identify all applications from System B that reference the "worst" device shared with System A.

Chapter 5.5: Analysis of VSAM data sets

CPEXpert can perform an analysis of problems with VSAM data sets, **if MXG has created the performance data base**. This section describes how this analysis is accomplished.

VSAM data set activity typically accounts for a large percent of I/O activity (more than 70% at some sites). Tuning of a few files or correcting common problems often can result in significantly improved performance (IBM benchmarks show up to 90% improvement resulting from some simple changes).

With CPEXpert Release 12.2, the DASD Component has been enhanced to provide a rudimentary analysis of common VSAM problems. Additional analysis will be added in future enhancements to the DASD Component.

The DASD Component can optionally analyze VSAM data set performance problems or potential problems under the following conditions:

- MXG has created the performance data base. Unfortunately, there is insufficient data retained in a NeuMICS performance data base to allow CPEXpert to analyze VSAM performance problems.
- The performance data base contains the MXG TYPE64 file (and CPEXpert is provided with the **%LET TYPE64=Y**; guidance variable).
- Most analysis of VSAM data sets also requires the MXG TYPE64 file (and CPEXpert is provided with the **%LET TYPE64=Y**; guidance variable). This is because the analysis depends on being able to identify the VSAM file type (e.g., KSDS, VRRDS, etc.), and the buffering technique used (e.g., NSR or LSR).

Section 2: Installing the DASD Component

Most of the DASD Component will be installed as normal part of installing CPExpert. However, up to three additional steps may be taken: (1) installing a modification to MXG or to NeuMICS to collect application-related or data set-related information, (2) defining workload categories, and (3) specifying data set information.

- **Collecting application-related or data set-related information.** This is an optional step. The DASD Component can optionally identify the specific applications (jobs and job steps) which access specific DASD volumes or use specific controllers or channel paths, or can analyze devices which provide a device I/O response time exceeding the response objective for selected data sets. This information can be quite useful in solving a DASD performance problem. With the application-related analysis, the DASD Component identifies the applications causing the problem or can identify the applications suffering from the problem. With the data set-related analysis, the DASD Component identifies devices which provide an I/O response time that exceeds a response objective for any critical data sets residing on the device. Providing this information is possible only if a modification is made to MXG or to MICS, to alter their processing of SMF Type 30(Data Definition) information.

The modification included with CPExpert is quite simple.

- For the application-related analysis, CPExpert records only the essential information required to identify the job name, the step name, the performance group number (for Compatibility Mode) or service class (for Goal Mode), the devices referenced and their SIO and connect times, and the beginning and end time of the measurement. This information is sufficiently concise that less than 25 cylinders of DASD are required to hold the information for all job steps executed in a relatively large installation. However, the information is sufficiently comprehensive for CPExpert to perform substantial analysis of DASD performance and the causes of poor DASD performance.
- For the data set-related analysis, CPExpert records information associated with specific DD names: the job name, the step name, the performance group number or service class, the device referenced and the SIO and connect time for the device, and the beginning and end time of the measurement.

Section 1 (Chapter 5) of this user manual describes the modifications and the processing of SMF information that CPExpert performs.

Chapter 1 describes how to install the modification for MXG code, and Chapter 2 describes how to install the modification for MICS code. The installation for either product is very straightforward.

- **Defining workload categories.** This is an optional step. CPExpert allows you to define categories of work, and have the analysis focus on specific workload categories.

This can be useful, for example, if you wish to identify your "loved ones" and have the analysis treat this workload category as more important than other work.

Chapter 3 of this section describes how to define workload categories to CPExpert. |

- **Specifying data set names.** This is an optional step. CPExpert allows you to specify data set names and associate a response objective with each data set. CPExpert will analyze SMF data to detect periods when devices on which the data sets reside provide an I/O response greater than the response objective. If any such intervals are detected, CPExpert will perform a detailed analysis of the devices to identify why the device I/O response exceeded the response objective.

Chapter 4 of this section describes how to define data set names and response objectives to CPExpert. |

The DASD Component provides many processing options (for example, you may exclude certain volumes from analysis, select particular time intervals for analysis, etc.). These processing options are not, strictly speaking, a part of the installation process. They may be employed as you refine your analysis of the performance of your DASD configuration. These processing options are described in Section 3 of this User Manual.

The instructions in this section assume that you have installed the CPExpert software. If you have not installed CPExpert, please install the software before continuing. (The procedures for installing CPExpert are described in the *CPExpert Installation Guide*.)

Chapter 1: Installing the modification for MXG

This chapter applies only if you use MXG to maintain your performance data base.

With normal MXG processing, the MXG TYPE30_D data set is not written because the data set would generally require **far** too much DASD space (often an entire volume would be required to hold the data set). Consequently, Barry Merrill does not write the TYPE30_D data set and he includes appropriate warnings in the IMAC30DD module that the data set may be very large. The CPExpert modification to MXG can allow you to collect essential information relating the DASD use to specific job steps and performance groups or service classes (with Goal Mode), without requiring massive amounts of DASD space. Additionally, the modification can allow you to collect information required to relate I/O activity to specific data sets.

Installing the modification for MXG is an optional step; however the DASD Component will be unable to relate DASD use and performance constraints to specific applications or to specific data sets unless you install the modification.

Implementing the DASD Component modification to MXG consists of five steps:

Step 1: Install the DASDMXG code.

Install the DASDMXG code into the SAS source library containing user modifications to MXG. This is done by copying the DASDMXG member from the CPEXPERT.SOURCE PDS into MXG.USERLIB (or whatever you call your library for user modifications to MXG). The DASDMXG member contains all the macros that will be invoked by the modifications described below.

Step 2: Modify MXG modules.

Modify one of two MXG modules to cause the DASD Component to output the DASD30DD data set. The MXG modules that may be modified are the IMACINTV module or the EXTY30U4 module.

- If you collect SMF Type 30 Interval records, modify the MXG module IMACINTV by including the following statement in the module:

```
%CPEOUTDD(V);    /* CPEXPERT MODIFICATION */
```

Exhibit 2-1 illustrates the MXG IMACINTV module with the modification.

```

/*  COPYRIGHT (C) 1985,1992 BY MERRILL CONSULTANTS DALLAS TEXAS      */
/*                                                                    */
/*  IMACINTV CONTROLS THE CREATION OF TYPE30_V DATASET, BUILT FROM    */
/*  TYPE 30 INTERVAL RECORDS (SUBTYPES 2 AND 3).                      */
/*  NOTE THAT THE PDB.SMFINTRV DATASET IS A RENAME OF TYPE30_V AND    */
/*  THUS CHANGING THIS MEMBER TO CREATE TYPE30_V OBSERVATIONS WILL    */
/*  ALSO CREATE PDB.SMFINTRV IF YOU EXECUTE BUILDpdb/BUILDpd3.        */
/*  BY DEFAULT, THIS MEMBER CREATES NO OBSERVATIONS IN TYPE30_V.      */
/*  IF YOU WANT ANY OBSERVATIONS TO BE CREATED IN THE DATASET, YOU    */
/*  MUST EXECUTE AN OUTPUT _LTY30UV; STATEMENT IN THIS MEMBER.        */
/*  FOR EXAMPLE, IF YOU WANT AN OBSERVATION IN TYPE30_V FOR EACH AND  */
/*  EVERY TYPE 30 SUBTYPE 2 OR SUBTYPE 3 RECORD READ, SIMPLY REMOVE  */
/*  THE COMMENT BLOCK AROUND THE FOLLOWING STATEMENT:                 */
/*  FOLLOWING STATEMENT:                                              */
/*                                                                    */
/*          OUTPUT _LTY30UV ;                                         */
/*                                                                    */
/*  THE _LTY30UV MACRO, DEFINED IN IMAC30, DEFAULTS TO TYPE30_V.      */
/*                                                                    */
/*  ALTERNATIVELY, IF YOU WISH TO CREATE OBSERVATIONS ONLY FOR       */
/*  STARTED TASKS, REMOVE THE COMMENT BLOCK FROM THE FOLLOWING       */
/*  STATEMENT:                                                        */
/*                                                                    */
/*          IF TYPETASK='STC' THEN OUTPUT _LTY30UV;                  */
/*                                                                    */
/*  IMACINTV IS %INCLUDED AFTER THE "IMPORTANT" VARIABLES, SUCH AS   */
/*  JOB, READTIME, STEPNAME, PROGRAM, TYPETASK, JOBCLASS, SMFTIME,    */
/*  ETC. HAVE BEEN DEFINED. THUS YOUR SELECTION CRITERIA CAN CAUSE   */
/*  TYPE30_V TO CONTAIN OBSERVATIONS ONLY FROM SPECIFIC JOBS AND/OR  */
/*  PARTICULAR PROGRAM NAMES, SUCH AS                                */
/*                                                                    */
/*          IF JOB='CICS' OR PROGRAM='DFHSIP' THEN OUTPUT _LTY30UV;   */
/*                                                                    */
/*  REALIZE THAT THIS MEMBER CONTROLS ONLY THE SELECTION OF TYPE 30  */
/*  INTERVAL RECORDS. YOU MUST (IN PARMLIB) SPECIFY THAT YOU DESIRE  */
/*  INTERVAL RECORDS TO BE CREATED BY SMF; YOU CAN SPECIFY SEPARATE  */
/*  DURATIONS INDIVIDUALLY FOR JOBS (JOB), TSO SESSIONS (TSU), AND   */
/*  OR STARTED TASKS (STC).                                          */
/*                                                                    */
/*                                                                    */
/*                                                                    */
/*                                                                    */
/* ***** THIS MODULE IS INCLUDED BY THE FOLLOWING MXG MODULES: ***** */
/* ***** VMAC30 ***** */
/* ***** */
/* ***** */
%CPEOUTDD(V);          /* CPEXPERT MODIFICATION */

```

MXG IMACINTV MODULE, WITH CPEXPERT MODIFICATION

EXHIBIT 2-1

MXG coding printed with permission of Merrill Consultants, Dallas, Texas.

The macro parameter "(V)" in the above code causes SAS statements to be generated to identify the source of the data (namely, the SMF Type 30 Interval records).

Collecting interval data is BY FAR the preferred approach, since much more comprehensive analysis can be accomplished.

- If you do **not** collect interval data but collect only Workload Termination data (that is, you collect SMF Type 30 Workload Termination records but do not collect Type 30 Interval records), you can still obtain information about the DASD use of specific job steps. However the analysis will not be as comprehensive as if you collect interval information.

If you collect only Workload Termination records, modify the MXG module EXTY30U4 by including the following statement in the module:

```
%CPEOUTDD(4); /* CPEXPERT MODIFICATION */
```

The macro parameter "(4)" in the above code causes SAS statements to be generated to identify the source of the data (namely, the SMF Type 30 Workload Termination records).

Exhibit 2-2 illustrates the MXG EXTY30U4 module with the modification.

```
/* COPYRIGHT (C) 1985,1992 BY MERRILL CONSULTANTS DALLAS TEXAS USA */
/*****
/* MEMBER EXTY30U4 - MXG OUTPUT EXIT FOR DATA SET TYPE30_4 */
/*****
/* THIS MEMBER CONTAINS THE OUTPUT STATEMENT FOR THE DATA SET. */
/* CHANGES IN DDNAME, DSNNAME, AND THE CREATION OF NEW VARIABLES */
/* CAN BE CODED HEREIN. */
/*****
OUTPUT _LTY30U4; /* TYPE30_4, DEFINED IN IMAC30 */

%CPEOUTDD(4); /* CPEXPERT MODIFICATION */
```

MXG EXTY30U4 MODULE, WITH CPEXPERT MODIFICATION

EXHIBIT 2-2

MXG coding printed with permission of Merrill Consultants, Dallas, Texas.

Step 3: Modify MXG IMAC30DD module.

Modify the MXG IMAC30DD module to cause the DASD Component to collect the Data Definition information that is contained in the DASD30DD data set. This is accomplished by including the following statement in the module:

```
%DASD_MXG; /* CPEXPERT MODIFICATION */
```

Exhibit 2-3 illustrates the MXG IMAC30DD module with the modification.

```

/* COPYRIGHT (C) 1985,1992 BY MERRILL CONSULTANTS DALLAS TEXAS */
/*****MEMBER=IMAC30DD*****/
/*
/* THIS MODULE CONTROLS THE EXISTENCE OF DATA SET   TYPE30_D
/* (WHICH IS RENAMED PDB.DDSTATS BY BUILDPDB/BUILDPEDE).
/*
/* BECAUSE TYPE30_D COULD CONTAIN ONE OBSERVATION FOR EVERY DD
/* CARD IN EVERY SUBTYPE OF EVERY TYPE 30 RECORD, THE DATASET
/* COULD BE VERY, VERY LARGE, AND THUS BY DEFAULT, MXG DOES NOT
/* CREATE ANY OBSERVATIONS IN TYPE30_D.
/*
/* YOU MUST MODIFY THIS INSTALLATION MACRO TO CREATE OBSERVATIONS!
/*
/* NOTE THAT DATASET TYPE30_D AND PDB.DDSTATS WILL ALWAYS EXIST;
/* THIS EXIT CONTROLS ONLY IF THERE ARE OBSERVATIONS CREATED.
/*
/* THE DSNNAME MACRO _LTY30UD IS DEFINED IN IMAC30 AND DEFAULTS
/* TO TYPE30_D.
/*
/* THE FOLLOWING EXAMPLE HAS BEEN COMMENTED OUT SO THAT IT IS NOT
/* EXECUTED, AND IT SHOWS ONE OF MANY CONDITIONS YOU MIGHT CHOOSE
/* TO RESTRICT WHICH DD SEGMENT'S OBSERVATIONS ARE OUTPUT.  THE
/* EXAMPLE WILL CREATE OBSERVATIONS ONLY FOR THE STEP TERMINATION
/* (SUBTYPE=4) RECORD, AND ONLY FOR THE SPECIFIC DEVICE NUMBER
/* (UCB ADDRESS=38F) IN THIS EXAMPLE.  YOU COULD CHOOSE TO TEST
/* FOR JOB, STEP, DDNAME, OR ANY OF THE VARIABLES IN TYPE30.
/*
/* NOTE: THE EXAMPLE TEST IS INSIDE A PAIR OF COMMENT SYMBOLS
/* WHICH PRECEED AND FOLLOW THE EXAMPLE "IF STATEMENT".
/* YOU MUST REMOVE THE COMMENT SYMBOLS AROUND THE TEST, AND
/* CHANGE THE TEST AS DESIRED, TO CREATE OBSERVATIONS.
/*
/*
/* IF DEVNR=038FX AND SUBTYPE=4 THEN OUTPUT _LTY30UD;
/*
/*****
/**** THIS MODULE IS INCLUDED BY THE FOLLOWING MXG MODULES: ****/
/****
/**** VMAC30 ****/
/****
/*****
/*
%DAED MXG: * CPEXPRT MODIFICATION TO ACQUIRE DAED INFORMATION:

```

MXG IMAC30DD MODULE, WITH CPEXPERT MODIFICATION

EXHIBIT 2-3

MXG coding printed with permission of Merrill Consultants, Dallas, Texas.

Step 4(alternate): Modify MXG EXPDBINC module and EXPDBVAR module.

This step applies if you use the standard MXG BUILDpdb or use SAS/ITSV to process SMF data.

Modify the MXG EXPDBINC module as shown below:

```
%INCLUDE SOURCLIB(DASDMXG); /* CPExpert modification */
```

Modify the MXG EXPDBVAR module as shown below:

```
MACRO _VARUSER /* CPExpert modification */  
  %%VARDDCPE(MAXDASD=nn,MAXDD=nnn,MINIO=nnn,COLLECT=xxxxxxx)  
%
```

The %%VARDDCPE macro statement generates code to define the DASD30DD data set, define the DASD30DS data set, and describe the variables contained in the KEEP statement for the data sets. (**NOTE:** two "%" symbols are required since the code is placed into SAS "old style" macros.) The macro parameters specify the maximum number of unique DASD devices referenced by any job step in your environment, the maximum number of DD statements associated with critical data sets, the significant number of I/O operations directed to any data set in any Type 30(Interval) record, and which type of data collection should be performed.

- The **MAXDASD** macro parameter specifies the maximum number of unique DASD devices that are referenced by any job step in your environment. Note that this is **not** the number of DD statements. Rather, this is the number of unique devices. There often will be many DD statements for a single device, but the DASD_MXG software combines these DD statements into one observation in the DASD30DD data set.

The default specification for the MAXDASD macro parameter is **MAXDASD=25**. The DASD_MXG software will produce a warning message if this value is too low, and the software simply ignores devices for any job step exceeding the value specified. (There is no other effect; that is, the software will not abort and the analysis software will continue to function properly.)

The only reason for specifying 25 versus 200 (for example) as an initial value, is that the larger the value, the more DASD space that will be used for the DASD30DD data set. The record size will increase by 15 bytes for each unique device specified. If you process 20,000 job steps each day and specify a maximum of 25 unique devices per job step, the total temporary DASD space required for the DASD30DD data set will be about 25 cylinders (3380). This is not a large amount in any event, but there is no point in writing larger records than are needed.

Therefore, MAXDASD=25 should be specified initially and the value should be increased only if required.

- The **MAXDD** macro parameter specifies the maximum number of unique DD names that may be associated with critical data sets. As described in Section 1 (Chapter 5.3), the CPExpert modification to MXG creates an array which contains the DD names of all DD statements associated with data sets which have been defined to CPExpert with a response objective. (The array is created only if you specify DDNAMES or BOTH in the COLLECT macro parameter described below.)

Each array element requires only 8 bytes of storage. This amount storage is insignificant. However, SAS generates a SAS data element name for each array element. SAS 5.18 has a restriction on the total number of data element names which can be defined. When this number is exceeded, SAS aborts with an error in the compilation stage. Therefore, you cannot set an extremely large value for the MAXDD parameter when executing under SAS 5.18.

The default specification for the MAXDD macro parameter is **MAXDD=200**, generating an array capable of containing 200 unique DD names. The DASD_MXG software will produce a warning message if the number of unique DD names is greater than the array size. The software will list all DD names which it could not process, and it simply ignores DD names for any job step exceeding the value specified. (There is no other effect; that is, the software will not abort and the analysis software will continue to function properly.)

If you receive a warning message from the CPExpert modification, you should consider (1) increasing the value of the MAXDD macro parameter or (2) restricting the number of data sets you define as critical.

- The **MINIO** macro parameter specifies the minimum number of I/O operations to any data set within a particular SMF Type 30(Interval) record. The CPExpert modification will ignore any data sets encountered if the number of I/O operations is below the MINIO value. The point of this parameter is that you probably will not worry about data sets with a very low activity, regardless of the DASD response time they experienced. There is no point in collecting information about very low activity data sets.

The default specification for the MINIO macro parameter is **MINIO=100**.

- The **COLLECT** macro parameter specifies whether to collect device information (required to associate device activity with application systems), to collect DD name information (required to associate device activity with data sets), or to collect both device and DD name information.

Specify **COLLECT=DEVICES** to collect device information. This is the default setting. If you have previously installed the modification for MXG, you do not have to make any changes to your installation.

Specify **COLLECT=DDNAMES** to collect DD name information.

Specify **COLLECT=BOTH** to collect both device information and DD name information.

WARNING: Please review Chapter 6 of this section before specifying DDNAMES or BOTH as an option for the COLLECT parameter. You should **not** specify either of these options unless you have defined critical data sets and have followed the procedures in Chapter 6.

Step 4(alternate): Modify the SAS job stream used to execute MXG.

This step applies if you use in-stream JCL for MXG to process SMF data. Modify the SAS job stream used to execute MXG by making the following simple modifications:

- Modify the %INCLUDE statement to add **DASDMXG** to the list of included code.
- Add the following macro statement after the MXG DATA statement:

%VARDDCPE(MAXDASD=nn,MAXDD=nnn,MINIO=nnn,COLLECT=xxxxxxx)

Please refer to the discussion in "Step 4 (if you use the standard MXG BUILDpdb to process SMF data)" for a description of the macro variables.

WARNING: Do not include a semicolon after the macro invocation!

Exhibit 2-4 illustrates a sample MXG SAS job stream with the three modifications highlighted. Of course, your own MXG SAS job stream may appear somewhat different than that shown in Exhibit 2-4. However, your MXG SAS job stream should be sufficiently similar to Exhibit 2-4 that you can appreciate how the modifications should appear in your own MXG SAS job stream.

Step 5: Add CPEDASD DD statement to the JCL.

Add the CPEDASD DD statement to the job control language (JCL) used to execute MXG. The CPEDASD DD statement defines the SAS library in which the CPExpert modification will place the DASD30DD data set or DASD30DS data set created by the modification.

The CPEDASD DD statement can refer to the CPEDATA SAS library if you wish, or you can allocate a different SAS library. For space allocation purposes, the CPEDASD data set requires slightly less than one cylinder of DASD space per 2,200 SMF Type 30(Interval) records processed by MXG. For example, if you typically have about 22,000 Type 30(Interval) records in your daily SMF file, you will need about 10 cylinders of DASD allocated to the CPEDASD SAS library. It is difficult to estimate the amount of space required for the DASD30DS data set, since this is a function of the number of critical data sets you define and the number of unique DD names associated with the data sets. This data set contains about 730 observations per track.

Exhibit 2-4 illustrates a sample MXG SAS job stream with the CPEDASD DD statement included, using the CPEDATA SAS library space.

```
//CPEDASD   DD   DSN=prefix.CPEXPRT.CPEDATA,DISP=OLD
//SYSIN     DD   *

%INCLUDE
SOURCLIB(VMACSMF,VMAC7072,VMAC71,VMAC73,VMAC74,VMAC75,VMAC77,
          VMAC78,VMAC30,VMACTSOM,VMAC434,VMAC40 VMAC434D,
          DASDMXG);
DATA

  %VARDDCPE(MAXDASD=25,MAXDD=200,MINIO=100,COLLECT=DEVICES)

  _VAR7072
  _VAR71

etc.

;
```

SAMPLE MXG SAS JOB STREAM

EXHIBIT 2-4

Chapter 2: Installing the modification for NeuMICS

This chapter applies only if you use MICS to maintain your performance data base.

The normal MICS processing summarizes the device information from the SMF Type 30(DD) EXCP segment into the BATWDA dataset. The BATWDA dataset is created only if specified during the MICS installation process. The BATWDA data set often is not created because it is so large, and the information is often too general for detailed performance analysis. CPEXpert does not use the MICS BATWDA file. Rather, CPEXpert uses data created by a proprietary modification to MICS. The CPEXpert modification to MICS is described in this section.

The CPEXpert modification to MICS allows you to collect essential information relating the DASD use to specific job steps and performance groups or service classes (with Goal Mode), without requiring massive amounts of DASD space. Additionally, the modification allows you to collect information required to relate I/O activity to specific data sets. The SAS data sets created by the modification to MICS normally are small (typically about 10 cylinders of DASD), and they are easy to process.

The CPEXpert modification to MICS invokes two standard MICS user exits and includes a simple (one-line) modification to the DYSMFFMT module. Installing the modification for MICS is an optional step; however the DASD Component will be unable to relate DASD use and performance constraints to specific applications or to specific data sets unless you install the modification.

Implementing the CPEXpert modification to MICS consists of five steps:

Step 1: Install the DASDMIC code.

Install the DASDMIC code into the prefix.MICS.USER.SOURCE library containing user modifications to MICS. This is done by copying the DASDMIC member from the CPEXPERT.SOURCE partitioned data set into the prefix.MICS.USER.SOURCE library. The DASDMIC member contains all the macros that will be invoked by the modifications described below.

Step 2: Modify the sharedprefix.MICS.USER.SOURCE(#SMFEXIT).

Modify the sharedprefix.MICS.USER.SOURCE(#SMFEXIT), by including SAS statements (1) to include the CPEXpert macros, (2) to define the _USRDMAP macro, and (3) to define the _USRSSF macro.

Exhibit 2-5 illustrates the normal sharedprefix.MICS.USER.SOURCE(#SMFEXIT) member before the modifications. Exhibit 2-6 illustrates the normal sharedprefix.MICS.USER.SOURCE(#SMFEXIT) member after the modifications.

```
/* some comments */  
%INCLUDE SOURCE(#SMFEXIT);
```

NORMAL sharedprefix.MICS.USER.SOURCE(#SMFEXIT) BEFORE MODIFICATION

EXHIBIT 2-5

```
/* some comments */  
%INCLUDE SOURCE(#SMFEXIT);  
  
/* CPEXPERT MODIFICATION */  
  
%INCLUDE USOURCE(DASDMIC);  
  
MACRO _USRDMAP  
/* DEVICE ADDRESS MAPPING EXIT */  
%%DASD_MIC;  
%  
  
MACRO _USSSFS  
/* STEP TERMINATION INTERIM SMF FILE EXIT */  
%%CPEOUTDD;  
%
```

NORMAL sharedprefix.MICS.USER.SOURCE(#SMFEXIT) AFTER MODIFICATION

EXHIBIT 2-6

WARNING: Please note that two (2) "%" characters are used when invoking "new style" macros inside "old style" macros in the user exits.

- The **%INCLUDE USOURCE(DASDMIC);** statement shown in Exhibit 2-6 causes MICS to include the DASDMIC macros as MICS loads its own code.
- The **%%DASD_MIC;** SAS macro statement shown as part of the _USRDMAP macro in Exhibit 2-6 generates code and arrays to process and store information related to job step use of DASD, by device address.
- The **%%CPEOUTDD;** SAS macro statement shown as part of the _USRSSFS macro in Exhibit 2-6 generates code to output the CPEDASD.DASD30DD data set, and clears the arrays described above.

The effect of Step 2 is that the normal _USRDMAP and _USRSSFS user exits included as null members in sharedprefix.MICS.SOURCE(#SMFEXIT) will be overridden by the code required for CPEExpert.

Step 3: Modify prefix.MICS.SOURCE(DYSMFFMT).

Modify the sharedprefix.MICS.SOURCE(DYSMFFM1) module¹ by including the following SAS statement just after the "DATA" statement:

```
%VARDDCPE(MAXDASD=nn,MAXDD=nnn,MINIO=nnn,COLLECT=xxxxxxx)
```

WARNING: Do not include a semicolon after the macro invocation!

The %VARDDCPE macro statement generates code to define the DASD30DD data set, define the DASD30DS data set, and describe the variables contained in the KEEP statement for the data sets. The macro parameters specify the maximum number of unique DASD devices referenced by any job step in your environment, the maximum number of DD statements associated with critical data sets, the significant number of I/O operations directed to any data set in any Type 30(Interval) record, and which type of data collection should be performed.

- The **MAXDASD** macro parameter specifies the maximum number of unique DASD devices that are referenced by any job step in your environment. Note

¹It is unfortunate that a modification must be made to the MICS coding, since this is contrary to standards implemented at many installations. A request was made to LEGENT Corporation (through the MICS Advisory Council, as Incident #441496) that a standard user exit be implemented after the DATA statement. If this request had been accepted and implemented, output files could be generated easily within the standard MICS user exit facility. Using a standard MICS user exit facility would obviate the need for a modification to MICS code. Computer Associates has not indicated a desire to make the modification.

that this is **not** the number of DD statements. Rather, this is the number of unique devices.

There often will be many DD statements for a single device, but the DASD_MIC software combines these DD statements into one observation in the DASD30DD data set.

The default specification for the MAXDASD macro parameter is **MAXDASD=25**. The DASD_MIC software will produce a warning message if this value is too low, and the software simply ignores devices for any job step exceeding the value specified. (There is no other effect; that is, the software will not abort and the analysis software will continue to function properly.)

The only reason for specifying 25 versus 200 (for example) as an initial value, is that the larger the value, the more DASD space that will be used for the DASD30DD data set. The record size will increase by 15 bytes for each unique device specified. If you process 20,000 job steps each day and specify a maximum of 25 unique devices per job step, the total temporary DASD space required for the DASD30DD data set will be about 25 cylinders (3380). This is not a large amount in any event, but there is no point in writing larger records than are needed. Therefore, MAXDASD=25 should be specified initially and this value should be increased only if required.

- The **MAXDD** macro parameter specifies the maximum number of unique DD names that may be associated with critical data sets. As described in Section 1 (Chapter 5.3), the CPExpert modification to MICS creates an array which contains the DD names of all DD statements associated with data sets which have been defined to CPExpert with a response objective. (The array is created only if you specify DDNAMES or BOTH in the COLLECT macro parameter described below.)

Each array element requires only 8 bytes of storage. This amount storage is insignificant. However, SAS generates a SAS data element name for each array element. SAS 5.18 has a restriction on the total number of data element names which can be defined. When this number is exceeded, SAS aborts with an error in the compilation stage. Therefore, you cannot set an extremely large value for the MAXDD parameter when executing under SAS 5.18.

The default specification for the MAXDD macro parameter is **MAXDD=200**, generating an array capable of containing 200 unique DD names. The DASD_MIC software will produce a warning message if the number of unique DD names is greater than the array size. The software will list all DD names which it could not process, and it simply ignores DD names for any job step exceeding the value specified. (There is no other effect; that is, the software will not abort and the analysis software will continue to function properly.) If you

receive a warning message from the CPExpert modification, you should consider (1) increasing the value of the MAXDD macro parameter or (2) restricting the number of data sets you define as critical.

- The **MINIO** macro parameter specifies the minimum number of I/O operations to any data set within a particular SMF Type 30(Interval) record. The CPExpert modification will ignore any data sets encountered if the number of I/O operations is below the MINIO value. The point of this parameter is that you probably will not worry about data sets with a very low activity, regardless of the DASD response time they experienced. There is no point in collecting information about very low activity data sets.
The default specification for the MINIO macro parameter is **MINIO=100**.

- The **COLLECT** macro parameter specifies whether to collect device information (required to associate device activity with application systems), to collect DD name information (required to associate device activity with data sets), or to collect both device and DD name information.

Specify **COLLECT=DEVICES** to collect device information. This is the default setting. If you have previously installed the modification for MICS, you do not have to make any changes to your installation.

Specify **COLLECT=DDNAMES** to collect DD name information.

Specify **COLLECT=BOTH** to collect both device information and DD name information.

WARNING: Please review Chapter 6 of this section before specifying DDNAMES or BOTH as an option for the COLLECT parameter. You should **not** specify either of these options unless you have defined critical data sets and have followed the procedures in Chapter 6.

Exhibit 2-7 illustrates the prefix.MICS.SOURCE(DYSMFFM1) code with the modification, reflecting explicit stating of the default values.

DATA

%VARDDCPE(MAXDASD=25,MAXDD=200,MINIO=100,COLLECT=DEVICES)

%MAINWRK(CCC=SMF,WRK=INI)

**SAMPLE sharedprefix.MICS.SOURCE(DYSMFFM1)
WITH CPEXPRT MODIFICATION**

EXHIBIT 2-7

Step 4: Add the CPEDASD DD statement to the JCL.

Add the CPEDASD DD statement to the job control language (JCL) used to execute MICS. The CPEDASD DD statement defines the SAS library in which the CPEExpert modification will place the DASD30DD and DASD30DS data sets created by the modification.

Exhibit 2-8 illustrates a sample CPEDASD DD statement.

//CPEDASD DD DSN=prefix.CPEXPRT.CPEDATA,DISP=OLD

SAMPLE CPEDASD DD STATEMENT IN MICS JCL

EXHIBIT 2-8

The CPEDASD DD statement can refer to the CPEDATA SAS library if you wish, or you can allocate a different SAS library. For space allocation purposes the CPEDASD data set requires slightly less than one cylinder of DASD space per 2,200 SMF Type 30(Interval) records processed by MXG or MICS to contain the application-related information written to the DASD30DD SAS data set (that is, if you have specified COLLECT=DEVICES or COLLECT=BOTH). For example, if you typically have about 22,000 Type 30(Interval) records in your daily SMF file, you will need about 10 cylinders of DASD allocated to the CPEDASD SAS library to contain the application-related information.

It is difficult to estimate the amount of space required for the data set-related information written to the DASD30DS data set, since this is a function of the number of critical data sets you define and the number of unique DD names associated with the data sets. The DASD30DS data set contains about 730 observations per track.

MICS users may wish to place the CPEDASD DD statement into PROCLIB and then do a MICS JCLGEN to copy the information into the MICS prefix.CNTL libname.

We STRONGLY suggest that you test the modification in a MICS Test environment before placing the modification into a MICS production environment! If you (or we) have made a mistake, the MICS production run might abort. This certainly would not generate a favorable impression about CPExpert.

Chapter 3: Defining workload categories

The Workload Definition Section of the DASGUIDE module can be used to (1) define workload categories, (2) associate performance groups with the workload categories prior to MVS (Goal Mode), (3) associate service classes with the workload category with MVS (Goal Mode), (4) assign a relative importance to the workloads, and (5) direct CPExpert to perform analysis based upon the defined workload categories.

The workload definition is optional; you do not have to define any workloads. However, you must define workloads if you wish CPExpert to analyze the DASD performance effects of contending workloads. Exhibit 2-9 illustrates a sample workload definition.

```
***** ;
*   WORKLOAD DEFINITION BY PERFORMANCE GROUP   ;
***** ;
* DO NOT REMOVE THE FOLLOWING SAS MACRO COMMENT LINE! ;
/* WORKLOAD DEFINITION
BEGIN THE WORKLOAD DEFINITION;
BATCH= 1           * SAMPLE: BATCH IS PGN 1           ;
TSO  = 2,6,9       * SAMPLE: TSO IS PGN 2,6,9         ;
CICS = 12,14       * SAMPLE: CICS IS PGN 12,14        ;
* DO NOT REMOVE THE FOLLOWING MACRO COMMENT LINE!    ;
*/
***** ;

***** ;
*   WORKLOAD DEFINITION BY SERVICE CLASS       ;
;
***** ;
* DO NOT REMOVE THE FOLLOWING SAS MACRO COMMENT LINE! ;
/* WORKLOAD DEFINITION (GOAL MODE) BEGIN THE WORKLOAD DEFINITION;
BATCH = BATHI,BATMED,BATLOW
CICS      = CICSAOR1,CICSAOR2,CICSAOR3
* DO NOT REMOVE THE FOLLOWING MACRO COMMENT LINE!    ;
*/
```

SAMPLE DISPLAY OF CPEXPRT.USOURCE(DASGUIDE) MODULE

EXHIBIT 2-9

he workload category definition syntax **prior** to MVS (Goal Mode) is:

NAME = PGN1[,PGN2...,PGNn]

- **NAME** is the workload name (e.g., BATCH, TSO, CICS, etc.). The name can be any number of characters (including blanks), since it is delimited by the equal sign. However, you should generally restrict the length to 10 characters, since some reports produced by CPExpert assume a maximum length for the column containing the workload name. No more than 10 workload category names may be defined. A workload name must not begin with "*".
- **PGN** is the number of the performance group(s) associated with this workload category. The performance group numbers are separated by commas, and the entire set may optionally be enclosed in parentheses. Up to 10 performance groups may be associated with a workload category². Any performance groups not associated with a workload category are placed in the OTHER category.

The workload category definition syntax **with** MVS (Goal Mode) is:

NAME = SRV1[,SRV2...,SRVn]

- **NAME** is the workload name (e.g., BATCH, TSO, CICS, etc.). The name can be any number of characters (including blanks), since it is delimited by the equal sign. However, you should generally restrict the length to 10 characters, since some reports produced by CPExpert assume a maximum length for the column containing the workload name. No more than 10 workload category names may be defined. A workload name must not begin with "*".
- **SRV** is the service class name(s) associated with this workload category. The service class names are separated by commas, and the entire set may optionally be enclosed in parentheses. Up to 10 service classes may be associated with a workload category³. Any service classes not associated with a workload category are placed in the OTHER category.

NOTE: If you are running Goal Mode, please be sure that you have modified USOURCE(GENGUIDE) to specify **%LET GOALMODE=Y**; so the CPExpert code related to service classes will be generated.

²The limits of 10 different workload categories and 10 performance groups per category were selected because they seem so large that they should satisfy your performance analysis needs. Please give us a call if you have a requirement for more than these limits.

³The limits of 10 different workload categories and 10 service classes per category were selected because they seem so large that they should satisfy your performance analysis needs. Please give us a call if you have a requirement for more than these limits.

Summary of Coding Restrictions:

- A workload name must **not** begin with "*".
- A maximum of 10 workload categories may be defined.
- A maximum of 10 performance groups or service classes may be assigned to any workload category.
- The entire definition for a workload category must be contained within a single line⁴.
- You may define the same performance groups or service classes in more than one workload category. However, the analysis may give strange results if the workload categories are defined as "competing" with each other (you have the strange situation of a performance group or service class competing with itself!).

There is a potentially significant limitation with defining workload categories: workload categories can be identified only by performance group or service class rather than performance group or service class **period**. This limitation is imposed by the SMF data; the standard SMF data associates performance or resource use to performance group periods or service class periods only in the SMF Type 72 records. It is not possible to associate DASD use, for example, to specific performance group or service class periods.

⁴This did not seem to be a particularly restrictive design constraint, so it didn't seem necessary to incorporate "continuation line" logic. If this does pose a problem, please give Computer Management Sciences a call.

Chapter 4: Defining critical data sets

The Data Set Definition Section of the DASGUIDE module can be used to (1) define data sets and (2) associate a response objective with each data set.

The data set definition is optional; you do not have to define any data sets. However, you must define data sets and specify a response objective for each data set if you wish CPEXpert to analyze the DASD performance provided to critical data sets.

Exhibit 2-10 illustrates a sample data set definition.

```
*****.
* DATA SET NAME/RESPONSE OBJECTIVE DEFINITIONS
;
*****.
* DO NOT REMOVE THE FOLLOWING SAS MACRO COMMENT LINE! ;
/* DEFINE DATA SET NAMES ;
DSN=USER.EXTERNAL.DATA,RESPONSE=20 ;
DSN=USER.INTERNAL,RESPONSE=30 ;
DSN=USER.INTERNAL.EXECUTIVE.DATA,RESPONSE=5;
DSN=USER.INTERNAL.SYSTEMS.DATA,RESPONSE=20;
* DO NOT REMOVE THE FOLLOWING MACRO COMMENT LINE! ;
*/
*****.
```

SAMPLE DISPLAY OF CPEXPERT.USOURCE(DASGUIDE) MODULE

EXHIBIT 2-10

The data set definition syntax is:

DSN = data.set.name,RESPONSE=obj

- **DSN** is the data set name. The data set name can be any valid data set name. The processing of data set name "wild cards" and the order in which data set names are processed are described below.
- **RESPONSE** is the device response objective for the identified data set. This specification is in milliseconds. For example, RESPONSE=20 means that the

device on which the data set resides should provide an average I/O response of 20 milliseconds during any RMF measurement interval. A response objective must be specified for each data set defined.

CPEXpert places each data set name into SAS macro variables, and places the associated response objective into macro variables corresponding to the data set name.

- CPEXpert processes the data set names with an implicit trailing "wild card" specification. That is, if only high level qualifiers are specified, all data set names beginning with the high level qualifiers are considered to be associated with the specified response objective.
- The order of appearance in the list controls the application of response objectives.

Examine Exhibit 2-10 to appreciate the significance of the above rules.

- All data sets beginning with "USER.EXTERNAL.DATA" will have a response objective of 20 milliseconds.
- All data sets beginning with "USER.INTERNAL" will have a response objective of 30 milliseconds. However, any data set beginning with "USER.INTERNAL.EXECUTIVE.DATA" will have a response objective of 5 milliseconds, while any data set beginning with "USER.INTERNAL.SYSTEMS.DATA" will have a response objective of 20 milliseconds.
- If the statement "DSN=USER.INTERNAL,RESPONSE=30" in Exhibit 2-10 were to be placed at the end of the list, it would override the previous statements and all data sets beginning with "USER.INTERNAL" would have a response objective of 30 milliseconds, regardless of whether they were associated with "EXECUTIVE.DATA" or with "SYSTEMS.DATA".

WARNING: Do not specify response objectives for data sets unless they are critical to accomplishing your management objectives, and do not make excessive use of the "wild card" technique. The purpose of this warning is to alert you to the fact that significant processing time and DASD space could be required if you do not judiciously apply this analysis technique.

The following operational notes apply to defining critical data sets:

- The DAS1415 module of CPEXpert must be executed before you execute MXG or MICS to perform your normal daily update of your performance data base. The DAS1415 module processes the data set name and response objectives you have defined in USOURCE(DASGUIDE). The DAS1415 module then processes the raw SMF Type 14/15 records to identify jobs which reference the data sets, and to build a SAS data set containing the DD names used to reference the data sets. Section 4 describes how to execute the DAS1415 module.
- You **cannot** add data set names to USOURCE(DASGUIDE) after you have executed MXG or MICS for any particular day, because these data set names will not have been processed by the DAS1415 module. Since the data set names were not processed by the DAS1415 module, the CPEXpert modification to MXG or MICS would not select data set information while MXG or MICS updated your performance data base.
- You **can** delete data set names or change the response objectives for data set names after your performance data base has been updated. For example, you may wish to do this between successive runs of the DASD Component of CPEXpert. The DASD Component will analyze DASD performance based upon current information in USOURCE(DASGUIDE). This is possible because the DASD Component processes USOURCE(DASGUIDE) to acquire current information. This information (data set names and response objectives) is used to guide the analysis of your performance data base.
- The VOLSER exclude or select logic described in Section 3 (Chapter 2.3 and Chapter 2.4, respectively) can be used to exclude or select specific volumes, without regard to the data set name processing.

Chapter 5: Defining Multiple PDBs

The Multiple Performance Data Base (PDB) section of the DASGUIDE module can be used to define multiple PDBs to CPEXpert.

The multiple PDB definition is optional; you do not have use or define multiple PDBs. This feature is applicable only to users who wish to use the DASD Component to analyze shared DASD contention problems (see Section 2, Chapter 2.5) **and** you have multiple performance data bases that contain information relating to systems that share DASD.

Exhibit 2-11 illustrates a sample definition of multiple performance data bases.

```
*****.
;
* MULTIPLE PERFORMANCE DATA BASE (PDB) DEFINITIONS
;
*****.
;
%LET MULTIPDB=Y;
%LET PDBLIB2 = PDBLIBB;
%LET PDBLIB3 = PDBLIBC;
*****.
;
```

SAMPLE DISPLAY OF CPEXPERT.USOURCE(DASGUIDE) MODULE

EXHIBIT 2-11

You should specify **%LET MULTIPDB=Y;** in prefix.CPEXPERT.USOURCE(GENGUIDE) to tell CPEXpert to process multiple performance data bases. You may then specify up to eight performance data bases in addition to the standard performance data base described by the PDBLIB DD statement, for a total of nine performance data bases.

You identify the DD statements that describe the additional performance data bases by specifying **%LET PDBLIBn = ddname;**, where "n" ranges from 2 through 9. In the example shown in Exhibit 2-11, the second performance data base is identified by the DD name of PDBLIBB, and the third performance data base is identified by the DD name of PDBLIBC. You can use any valid DD name to describe the performance data bases. The definitions must be in numerical order (that is, the first additional PDB must be described by PDBLIB2, the second additional PDB must be described by PDBLIB3, etc.).

Section 3: Specifying Guidance Variables

The CPEXPERT.USOURCE(DASGUIDE) partitioned data set member contains variables to establish the overall guidance for the DASD Component. You modify the variables in the DASGUIDE member whenever you wish to change the guidance to CPEXpert. This chapter describes these variables, how the variables are used, and what should be specified for the variables.

The variables in the DASGUIDE module can be viewed as "data selection and presentation" variables and "analysis control" variables. These two types of control variables are discussed separately.

The data selection and presentation variables allow you to select particular time intervals to be analyzed, and allow you to specify how the results from the analysis are to be presented.

The analysis control variables allow you to control the analysis the DASD Component will perform. These variables specify whether expanded analysis is to be performed, identify critical workloads, specify critical data sets, exclude devices from analysis, control CPEXpert's analysis of VSAM data set problems or potential problems, etc.

Please do not allow CPEXpert to perform analysis or produce reports that are meaningless in your environment. If the analysis and reports produced by CPEXpert do not meet your needs, alter the guidance to CPEXpert. If the guidance is insufficient, please call Computer Management Sciences at (703) 922-7027 (or e-mail Don_Deese@cpexpert.com) so we can make changes to improve CPEXpert for you!

Chapter 1: Data Selection and Presentation Variables

The data selection and presentation variables allow you to select particular time intervals to be analyzed, and allow you to specify how the results from the analysis are to be presented. This chapter describes these variables.

Exhibit 3-1 illustrates the data selection and presentation variables contained in USOURCE(DASGUIDE).

```
*****;
*   DATA SELECTION AND PRESENTATION VARIABLES   ;
*****;

%LET CONFIG    =&CPEWORK;      * SAS LIBRARY FOR I/O CONFIGURATION DATA ;
%LET CONFIGX   =;              * CONFIGURATION DATA ALREADY EXISTS   ;
%LET DASDATES  =01FEB1991;     * START DATE FOR DATA ANALYSIS    ;
%LET DASTIMES  =08:00:00;      * START TIME FOR DATA ANALYSIS    ;
%LET DASDATEE  =31DEC9999;     * END DATE FOR DATA ANALYSIS      ;
%LET DASTIMEE  =16:00:00;      * END TIME FOR DATA ANALYSIS      ;
%LET DASDAT2E  =0;             * DEFAULT SECOND SELECTION DATE - END ;
%LET DASDAT2S  =0;             * DEFAULT SECOND SELECTION DATE - START ;
%LET DASTIM2E  =0;             * DEFAULT SECOND SELECTION TIME - END ;
%LET DASTIM2S  =0;             * DEFAULT SECOND SELECTION TIME - START ;
%LET SHIFT     =Y;             * START AND END TIMES REFER TO A SHIFT ;
%LET SYSTEM    =*ALL;          * SPECIFY SYSTEM TO PROCESS (*ALL=ALL) ;
%LET SYSTEMn   =system;        * PROCESS SYSTEMn (n = 1-9)          ;
%LET SYSPLEX   =*ALL;          * SPECIFY SYSPLEX TO PROCESS (*ALL=ALL) ;
%LET SASODS    =N;             * CONTROLS WHETHER SAS ODS IS USED   ;
%LET PATH      =;              * PATH FOR ODS OUTPUT                ;
%LET FRAME     =DASFRAME;      * GENERIC ODS FRAME NAME            ;
%LET CONTENTS  =DASDCONT;      * GENERIC ODS CONTENTS NAME         ;
%LET BODY      =DASDBODY;      * GENERIC ODS BODY NAME             ;
%LET LINKPDF   =;              * LINK TO CPEXPRT DOCUMENTATION        ;
%LET STYLE     =;              * ODS HTML STYLE OPTION                ;
%LET PDFODS    =N;             * CONTROLS WHETHER SAS PDF IS USED   ;
%LET PDFFILE   = filename;     * DEFINES THE SAS PDF OUTPUT FILE    ;
%LET URL       =N;             * CONTROLS .HTM IN SAS ODS FRAME OUTPUT ;
%LET VERBOSE   =V;             * RESULTS: VERBOSE/CONCISE/SUMMARY  ;
%LET MAXRULES  =1000;          * PRODUCE -REPORT TOO BIG- GUIDANCE ;
```

DATA SELECTION AND PRESENTATION VARIABLES

EXHIBIT 3-1

Chapter 1.1: CONFIG variable

The CONFIG variable allows organizations to “save” the model of the I/O configuration that is created by CPExpert into a SAS library, described by the CONFIG DDNAME. The default value for the CONFIG guidance variable is **%LET CONFIG=WORK;** which causes the I/O configuration model to be placed into the SAS WORK temporary library. You can cause the DASD Component to place the I/O configuration into a permanent SAS library by using the CONFIG guidance variable. The I/O configuration model created by CPExpert contains not only the I/O configuration (physical channels, logical channels, control units, and devices attached to each system), but also contains key performance metrics related to each aspect of the I/O configuration. This option can be useful for reporting I/O activity for capacity planning. If you wish to save the I/O configuration model produced by CPExpert, specify **%LET CONFIG=ddname;**, where “ddname” is the name of JCL DD statement that describes the SAS library where the I/O configuration model is to be placed.

Chapter 1.2: CONFIGX variable

The CONFIGX variable allows organizations to specify that the I/O configuration model already exists in the library described by the CONFIG variable (see above). When **%LET CONFIG=EXIST;** is specified, CPExpert does not create the I/O configuration model, but assumes that a current model is present. This option eliminates the overhead associated with creating the I/O configuration model. Since the I/O configuration model contains measurement data acquired from the performance data base, you must be sure that the I/O configuration model represents measurement data that you wish CPExpert to analyze.

Chapter 1.3: DASDATES and DASTIMES variables

The DASDATES and DASTIMES variables specify the start date and start time, respectively, for the interval of SMF data the DASD Component is to analyze. These variables (in conjunction with the DASDATEE and DASTIMEE variables) allow you to select specific periods of data to analyze. The following example shows how to specify that data selection should start at 08:00:00 on March 6, 2001:

```
%LET DASDATES = 06MAR2001; * START DATE FOR DATA ANALYSIS;  
%LET DASTIMES = 08:00:00; * START TIME FOR DATA ANALYSIS;
```

Chapter 1.4: DASDATEE and DASTIMEE variables

The DASDATEE and DASTIMEE variables specify the end date and end time, respectively, for the interval of SMF data the DASD Component is to analyze. The following example shows how to specify that data selection should end at 17:00:00 on March 8, 2001:

```
%LET DASDATEE = 08MAR2001 * END DATE FOR DATA ANALYSIS;  
%LET DASTIMEE = 17:00:00; * END TIME FOR DATA ANALYSIS;
```

Chapter 1.5: DASDAT2S and DASTIM2S variables

The **DASDAT2S** and **DASTIM2S** variables are optional. These variables specify the start date and start time, respectively, for a second interval of SMF data the DASD Component is to analyze. These variables (in conjunction with the optional DASDAT2E and DASTIM2E variables) allow you to select a second period of data to analyze, in addition to the period specified by the DASDATES/DASTIMES and DASDATEE/DASTIMEE selection variables. The following example shows how to specify that a second period of data selection should start at 20:00:00 on March 6, 2001:

```
%LET DASDAT2S = 06MAR2001 * START DATE FOR DATA ANALYSIS;  
%LET DASTIM2S = 20:00:00; * START TIME FOR DATA ANALYSIS;
```

Chapter 1.6: DASDAT2E and DASTIM2E variables

The **DASDAT2E** and **DASTIM2E** variables are optional. These variables specify the end date and end time, respectively, for a second interval of SMF data the DAS Component is to analyze. These variables (in conjunction with the optional DASDAT2S and DASTIM2S variables) allow you to select a second period of data to analyze, in addition to the period specified by the DASDATES/DASTIMES and DASDATEE/DASTIMEE selection variables. The following example shows how to specify that a second period of data selection should end at 22:00:00 on March 6, 2001:

```
%LET DASDAT2E = 06MAR2001; * END DATE FOR DATA ANALYSIS;  
%LET DASTIM2E = 22:00:00; * END TIME FOR DATA ANALYSIS;
```

Chapter 1.7: MAXRULES variable

A very large amount (several hundred pages) of output can be produced during some executions of the DASD Component. Such a large amount of output is not easy to analyze, and tends to cause the report to be rejected by analysts or management. Fortunately, there are several options in the DASD Component that can be used to reduce the amount of output. New users of the DASD Component might not be aware of these options, and might not be favorably impressed with the DASD Component since they would perceive the DASD Component output to be useless. This is not a desirable result!

The MAXRULES variable is used to control whether CPExpert detects that “a large amount of output” was produced. CPExpert produces suggestions on how to reduce the output to a more manageable size when the number of rules (and supporting data) exceed the value of the MAXRULES guidance variable¹.

The default value of the MAXRULES guidance variable is 1000, which should be approximately 50 pages of output (basic rules have an average of 10 lines per description, while supporting rules have only one line for the data). You can alter the MAXRULES guidance variable and cause CPExpert to produce the suggestions for ways to reduce the output when either more or less rules (and supporting data) are produced. For example, to indicate that the suggestions on ways to reduce the output should be printed only when more than 2000 rules were produced, specify:

```
%LET MAXRULES=2000 ; * PRODUCE “REPORT TOO BIG” GUIDANCE;
```

Chapter 1.8: SHIFT variable

The SHIFT variable is used with the DASDATES, DASTIMES, DASDATEE, and DASTIMEE variables. The SHIFT variable allows you indicate how the time-selection variables should be used.

- If the SHIFT variable is "N", the time-selection will be based upon the absolute start and end dates/times specified. For example, if you wish CPExpert to process **all** data during a week, the start date and start time would be specified as the beginning of the week, and the end date and end time would be specified as the end of the week. You would specify "%LET SHIFT = N;" to process each 24-hour day. In the example shown above, data would be processed from 08:00:00 on 4 March until 17:00:00 on 8 March.

¹It would be more elegant to allow users to specify the maximum number of pages produced. Sadly, the “too much output” discussion is produced at the front of the output (so it will be the first thing a user sees upon examining the output), and the number of pages in the output is not known until the output is produced. Thus, the MAXRULES variable must be an estimate of the amount of output.

- If the SHIFT variable is "Y", the time-selection will be based upon the start and end dates, and the start and end times within each selected date. In the example shown above, perhaps you wished to process only the daily shift beginning at 08:00:00 and ending at 17:00:00. You would specify "%LET SHIFT = Y;" to process only the identified shift data, during the selected dates.

Chapter 1.9: SYSTEM variable

The SYSTEM variable is used to specify whether all systems in the performance data base should be evaluated, or to select a specific system identification to be evaluated.

Some users have data from multiple systems in their performance data base. For many of these users, or for users who have data for a single system represented in their performance data base, the default "**ALL" will be appropriate. No change of the SYSTEM variable would be required for these users.

However, some users who have data from multiple systems may wish to evaluate only a single system with the parameters specified in this member of DASGUIDE. For example, they might be temporarily interested in evaluating the performance of only an "important" system (such as a major production system) and not be interested in evaluating the performance of other systems with data in the performance data base. This evaluation can be accomplished by changing the SYSTEM variable to specify the system identification to be evaluated. For example, to specify that only data from SYS1 should be evaluated, specify:

```
%LET SYSTEM = SYS1 ; * PROCESS ONLY DATA FROM SYS1;
```

In another situation, a CPExpert user might wish to evaluate different systems with different DASGUIDE parameters. These different evaluations can be accomplished by different executions of the DASD Component. For each execution of the DASD Component, the USOURCE DD statement would be changed to reference different USOURCE libraries. Each USOURCE library would contain guidance members with appropriate guidance variables. The SYSTEM variable for each DASGUIDE guidance member would specify the system identification to which the guidance applied.

If you are using MICS as your performance data base, the SYSTEM variable refers to the MICS system identifier, rather than the SMF system identifier.

Chapter 1.10: SYSTEMn variable(s)

The SYSTEMn variable(s) are used to select multiple systems to be evaluated.

As described in the SYSTEM guidance variable discussion above, some sites have data from multiple systems in their performance data base. These sites can process data from all systems by specifying %LET SYSTEM=ALL; in USOURCE(CICGUIDE), or can select a specific system to process by specifying %LET SYSTEM=system; in USOURCE(CICGUIDE), where "system" is the system identification of the system to be processed.

Some sites have data from multiple systems in their performance data base and do not want to process all systems, but do wish to process more than one system. For example, some systems might be production systems and some might be test systems. For these sites, the **SYSTEMn** guidance variable can be used to select more than one specific system to analyze.

The SYSTEM guidance variable can be used to select data from only one system to analyze, and the SYSTEMn guidance variable(s) can be used to select up to 9 additional systems to analyze. For example, if you wish to analyze data from four systems (named SYSA, SYSB, SYSC, AND SYSX) in a single execution of the DASD Component, specify:

```
%LET SYSTEM = SYSA ; * PROCESS DATA FROM SYSA;  
%LET SYSTEM1 = SYSB ; * PROCESS DATA FROM SYSB;  
%LET SYSTEM2 = SYSC ; * PROCESS DATA FROM SYSC;  
%LET SYSTEM3 = SYSX ; * PROCESS DATA FROM SYSX;
```

Chapter 1.11: SYSPLEX variable

The SYSPLEX variable is used to specify whether data from each sysplex in the performance data base should be evaluated, or whether CPExpert should select data for a specific sysplex to be evaluated.

Some users have data from more than one sysplex in their performance data base. For many of these users, or for users who have data for a single sysplex represented in their performance data base, the default "*"ALL" will be appropriate. No change of the SYSPLEX variable would be required for these users.

However, some users who have data from more than one sysplex may wish to evaluate only a single sysplex with the parameters specified in this member of DASGUIDE. This evaluation can be accomplished by changing the SYSPLEX variable to specify the sysplex

to be evaluated. The following example shows how to specify that data only from PRODPLEX should be evaluated:

```
%LET SYSPLEX = PRODPLEX ; * PROCESS DATA ONLY FROM PRODPLEX;
```

Chapter 1.12: SAS Output Delivery System

Output from CPEXpert is created using Basic SAS statements. This Basic SAS output is designed for a standard SAS printer (line) format. With SAS Release 8, SAS users can use the SAS Output Delivery System to create output that is formatted in Hypertext Markup Language (HTML). This output can be browsed with Internet Explorer, Netscape, or any other browser that fully supports the HTML 3.2 tag set.

The CPEXpert WLM Component, DB2 Component, CICS Component, and DASD Component support the SAS ODS features.

Please reference the *CPEXpert Installation Guide* for more detailed information about using the SAS ODS feature of CPEXpert.

Chapter 1.13: VERBOSE variable

The VERBOSE variable provides a control on the amount of narrative that the DASD Component lists with each rule result. Some installations prefer to produce only concise findings, and evaluate the results when the findings are significant. Other installations wish to produce expanded findings, and evaluate the results on a daily basis. You can use the VERBOSE variable to control the amount of narrative, depending upon your preferences. The options with the VERBOSE variable are:

- V** = Print verbose comments related to each rule that was produced during the evaluation of DASD performance. The verbose comments describe the rule, provide key information associated with the rule, and may provide a specific reference related to the rule.
- C** = Print "concise" results about the DASD devices. Each "rule record" created by the DASD Component has a "level" associated with it. "Level 1" records are those records describing basic information about the DASD devices with the most potential for performance improvement. "Level 2" records are those records that analyze the causes of poor performance and suggest ways in

which improvements may be made. When the **C** verbose option is specified, only Level 1 information is printed.

- S** = Print “summary” results about the DASD devices. When the **S** verbose option is specified, only DAS000 (sysplex-wide intensity) and DAS050 (system intensity) information is printed.

Chapter 2: Analysis Control Variables

The analysis control variables allow you to control how the DASD Component analyzes your DASD configuration. These variables can be used to (1) specify whether expanded analysis is to be performed based upon workload categories, (2) specify whether expanded analysis is to be performed based upon data set names, (3) exclude volumes from analysis altogether, or (4) select specific volumes for analysis.

The analysis control variables are optional. You do not have to specify any analysis control variables unless you wish to alter the basic processing performed by the DASD. This chapter describes the analysis control variables associated with the DASD Component. Exhibit 3-2 illustrates the analysis control variables contained in USOURCE(DASGUIDE). Please note that these are basic analysis guidance variables. Additional guidance variables are listed in other chapters in this section.

```

*****;
*   ANALYSIS GUIDANCE VARIABLES                               ;
*****;
%LET ANALYZE    = 3;    * DEFAULT FOR NUMBER OF DEVICES TO ANALYZE      ;
%LET DASDEXCP   =100;   * MINIMUM NUMBER OF EXCP TO SELECT (DASD30DD)    ;
%LET DASDSN     = N;    * ANALYZE CRITICAL DATA SETS?                   ;
%LET EVALDASD   =Y;     * OPTION TO PRODUCE ONLY I/O CONFIGURATION        ;
%LET EXCLUDE    = N;    * DEFAULT FOR EXCLUDING VOLUMES                  ;
%LET GOALMODE   = Y;    * OPERATING IN GOAL MODE                          ;
%LET LIST42DS   =10;    * LIST ONLY TOP NUMBER OF DATA SETS             ;
%LET LISTALL    =Y;     * LIST DATA FOR ALL RMF INTERVALS               ;
%LET LOVED1     =N;     * PERFORM 'LOVED ONE' ANALYSIS?                   ;
%LET LOVEDALL   =N;     * PERFORM 'LOVED ONE' ANALYSIS ON ALL DEVICES?    ;
%LET MIN42PCT   =0.1;   * EXCLUDE REPORTING LOW-ACTIVITY DATA SETS      ;
%LET MINIORT    =0      * MINIMUM I/O RATE TO ANALYZE                    ;
%LET MINIOWT    =0;     * MINIMUM I/O INTENSITY (OR WEIGHT) TO ANALYZE   ;
%LET MINRESP    =0;     * MINIMUM I/O RESPONSE TO ANALYZE                ;
%LET REPORT     =11;    * NUMBER OF TOTAL VOLUMES TO REPORT              ;
%LET SELECT     = Y;    * DEFAULT FOR VOLUME SELECTION                   ;
%LET SHARED     = N;    * DEFAULT FOR PROCESSING SHARED SYSTEM DATA     ;
%LET TYPE30DD   = N;    * MXG TYPE30 MODIFICATION (Y = AVAILABLE)        ;
%LET TYPE42DS   = N;    * MXG TYPE42DS (Y = AVAILABLE)                   ;
*****;

```

ANALYSIS GUIDANCE VARIABLES

EXHIBIT 3-2

Chapter 2.1: Number of devices to analysis: ANALYZE variable

The optional ANALYZE variable is used to specify the number of devices that CPExpert should analyze in detail.

As explained in Section 1, most Data Base Administrators are able to address a limited number of DASD problems each day. The DASD Component was designed to identify the most serious DASD problems. The initial design identified the device with the most potential for improvement, and analyzed this device in detail.

Some users of CPExpert felt that the "identify the device with the most potential for improvement" design philosophy was unnecessarily restrictive. These users wanted to specify the number of devices that should be analyzed in detail. The ANALYZE variable meets this requirement.

You can direct CPEXpert to analyze up to 999 DASD devices in a single execution, by using the ANALYZE variable to specify the number of devices to analyze in detail. CPEXpert will provide a detailed analysis of all devices whose device response time exceeded the average device response time for their device type, up to the number specified by the ANALYZE variable. For example, if you wish CPEXpert to analyze the top 25 devices with the worst performance, specify **%LET ANALYZE=25;** in USOURCE(DASGUIDE).

You likely will receive an unwieldy amount of output if you specify a large number of devices using the ANALYZE variable. One way to limit the output is to specify **%LET VERBOSE=S** so that only device summary results will be printed.

Chapter 2.2: Exclude devices with low activity: DASDEXCP variable

The modification to MXG or NeuMICS includes a MINIO parameter that specifies the minimum number of I/O operations to any data set within a particular SMF Type 30(Interval) record. The CPEXpert modification will ignore any data sets encountered if the number of I/O operations is below the MINIO value. The default for the MINIO parameter is MINIO=100, indicating that data sets will be included in the DASD30DD data set generated by the CPEXpert modification if at least 100 EXCPs were directed to the data set.

This screening is based on data set, and is not based on the total I/O activity to the device. Consequently, some devices might have data represented in the DASD30DD data set (generated by the CPEXpert modification), but these devices might have very low activity. For example, a device could be represented in the DASD30DD with only 100 EXCPs per SMF Type 30 recording interval if only a single data set were referenced 100 times.

It normally is not interesting to consider devices with a very low activity, and many sites have a large number of such devices in the DASD30DD data set created by CPEXpert. To reduce processing time, the DASDEXCP guidance variable can be used to exclude such low-activity volumes.

The default value for the DASDEXCP guidance variable is 100 (the same as the MINIO variable in the CPEXpert modification) that creates the DASD30DD data set. If you wish CPEXpert to ignore devices with low activity, change to DASDEXCP guidance to a larger value. For example, if you wish CPEXpert to consider only those devices with more than 1000 EXCPs per SMF Type 30 recording interval, specify **%LET DASDEXCP=1000;** in USOURCE(DASGUIDE).

Chapter 2.3: Analyze using response objectives: DASDSN variable

The DASGUIDE **DASDSN** guidance variable specifies whether to perform expanded analysis based upon data set response objectives. Specify **%LET DASDSN = Y**; to tell CPEXpert to perform expanded analysis based upon data set response objectives.

Please refer to Section 2 (Chapter 4) for operational notes regarding changing your specification of critical data sets or changing response objectives associated with the data sets.

Please refer to Section 3 (Chapter 5) for a description of the procedures to follow if you wish to perform analysis based on data set response objectives.

Chapter 2.4: Produce only I/O configuration: EVALDASD variable

The EVALDASD guidance variable allows organizations to suppress the evaluation of DASD I/O problems. This option is for those users who wish to use the DASD Component simply to produce SAS data sets representing the I/O configuration (and associated data) for each system in the sysplex, and a SAS data set for the overall sysplex.

The DASD Component processes SMF data (in a performance data base) to create a single record for each RMF measurement interval, that contains information about the channel paths, controllers, and device for each volume. This information is interesting to some users for management reporting purposes, or for data base administrator reporting purposes. These users do not necessarily want to use the DASD Component to analyze DASD problems, since their sole use of the DASD Component is to have an easy-to-use tool that creates the I/O configuration information. This option can be useful for reporting I/O activity for capacity planning.

- If the EVALDASD guidance variable is "Y" (which is the default value), the DASD Component will create the model of the I/O configuration and evaluate DASD performance problems represented by this model.
- If the EVALDASD guidance variable is "N", the DASD Component will create the model of the I/O configuration and place the results in the SAS library described by the optional **CONFIG DDNAME** DD statement. The DASD Component will then issue an **ABORT RETURN 100**; statement within the SAS coding. This statement results in SAS aborting the SAS session, with a return code of 100.

Chapter 2.5: Excluding volumes from analysis: EXCLUDE variable

The EXCLUDE guidance variable is used to tell CPEXpert whether you wish to exclude volumes from analysis. The EXCLUDE variable acts as a "switch" to control whether CPEXpert processes the DASGUIDE member searching for volumes to exclude from analysis. The point of having a "switch" variable is that some installations may wish to regularly exclude volumes from analysis, but periodically analyze all volumes.

If the EXCLUDE variable is **N**, CPEXpert will not exclude any volumes from analysis. If the EXCLUDE variable is **Y**, CPEXpert will process the DASGUIDE member to identify all volumes to be excluded.

Please refer to Chapter 3 for more information on excluding volumes from analysis.

Chapter 2.6: Specifying data sets to list: LIST42DS variable

If SMF Type 42 (Data Set Statistics) information is available² in a MXG performance data base, the DASD Component will process MXG TYPE42DS to select data sets that were referenced during RMF measurement intervals in which the poorly performing devices exceeded the average performance. There can be many data sets referenced on the devices (hundreds or even thousands of data sets can be referenced). It is not helpful to have a large number of data sets listed. Consequently, CPEXpert uses the **LIST42DS** variable to limit the number of data sets listed in any RMF measurement interval, for a particular device. Only the number of data sets specified by the LIST42DS variable will be listed individually, and any remaining data sets will be summarized and listed on a single line.

The default value for the LIST42DS variable is 10, indicating that 10 data sets will be listed for each RMF interval in which a poorly-performing device had a performance problem.

You can alter this number of data sets listed by using the LIST42DS guidance variable. For example, if you wish CPEXpert to list only 5 data sets for each RMF interval in which a poorly-performing device had a performance problem, specify **%LET LIST42DS=5;** in USOURCE(DASGUIDE).

Chapter 2.7: Perform "loved one" analysis: LOVED1 variable

The LOVED1 guidance variable is used to control whether the DASD Component performs "loved one" analysis.

²**%LET TYPE42DS = Y;** must be specified in USOURCE(GENGUIDE) or USOURCE(DASGUIDE) to advise CPEXpert that TYPE42DS is available.

IBM speakers at professional conferences have long advocated the performance analysis principle of "Know your loved ones and always have someone to kick around." As described in Section 1 (Chapter 5.2) and Section 2 (Chapter 5), you can identify workload categories (such as TSO, CICS, BATCH, etc.) to the DASD Component. The DASD Component can then analyze your DASD configuration from the view of the workload categories, focusing on the performance provided to your "loved one" workload.

Analysis of critical (or "loved one") workload is possible only if you accomplish the following:

- You have installed the modification to MXG or MICS to collect application-related information as described in Section 2 of this User Manual. The DASD Component cannot perform expanded analysis based on workload categories unless it has information relating DASD use to specific performance groups or service classes, jobs, and job steps. This information is extracted from the SMF Type 30(DD) records by the modifications described in Chapter 1 and Chapter 2 of Section 2.
- You have defined workload categories to the DASD Component. At least one workload category (the most critical category) must be defined using the process described in Chapter 3 of Section 2.
- You identify the critical (or "loved one") workload to the DASD Component using the **LOVED1** guidance variable.
- You specify **%LET TYPE30DD = Y;** in the USOURCE(GENGUIDE) CPExpert general guidance module. This tells CPExpert that you have installed the modification to MXG or MICS to collect application-related information.

If you accomplish the above steps, the DASD Component will perform an expanded analysis of your DASD configuration from the perspective of the "loved one" workload category. That is, only devices and channel paths used by the critical workload will be analyzed. This analysis can be quite useful to determine which DASD devices are providing the worst performance from the perspective of the critical workload.

The **LOVED1** guidance variable specifies whether to perform expanded analysis of workload categories, and specifies the most critical (or "loved one") workload.

For example, if your critical workload is CICS production executing in Service Class CICSOR1, Service Class CICSOR2, and Service Class CICSOR3 you might define the CICS workload as illustrated in Exhibit 2-9 in Section 2. You could identify the CICS workload as your most critical workload as:

```
%LET LOVED1 = CICS ; * CICS IS THE CRITICAL WORKLOAD;
```

In this example, the DASD Component will analyze all DASD use by Service Class CICSOR1, CICSOR2, and CICSOR3. The DASD Component will isolate the devices providing the worst performance from the view of this CICS workload, identify likely causes of poor performance, identify the job and job steps in **other service classes** that **probably** cause performance to be poor, and suggest ways to correct the problems. A similar analysis would be performed for performance groups, if you are operating in Compatibility Mode.

A significant phrase was used in the above paragraph: "probably cause performance to be poor". Please review the discussion associated with appropriate rules in Appendix A to appreciate why the identification might not be precise, but might be only "probable". In brief summary, the identification is "probable" rather than "precise" because of synchronization problems between SMF Type 30 data and SMF Type 70(series) data, unless you have specified that the SMF/RMF recording be synchronized. This issue is discussed in some detail in Section 5 and in appropriate rules in Appendix A.

The DASD Component will not perform "loved one" analysis if the LOVED1 macro variable is null (this is the default specification).

The "loved one" analysis is incompatible with the data set response objective analysis³. That is, if you perform "loved one" analysis, you cannot perform an analysis of response objectives for critical data sets in the same execution of CPExpert. You could perform "loved one" analysis in one execution of the DASD Component and perform data set response objective analysis in another execution of the DASD Component, processing the same data in your performance data base.

Chapter 2.8: Analyze all devices referenced by "loved one" applications: LOVEDALL variable

With the normal "loved one" analysis, the DASD Component selects devices for detailed analysis only if the performance of these devices is worse than the average for their device type (see Section 1 for a discussion of this process). With the special case of devices referenced by "loved one" work, some users wish to examine *all* devices referenced by the work, regardless of whether the device performance was worse than average. This ability is provided by the LOVEDALL guidance variable. When **%LET LOVEDALL=Y;** is specified, the DASD Component will analyze the top "n" devices referenced by the "loved one" work, regardless of whether the performance of these devices was worse than average for their device type. The "n" in this case is as specified by the ANALYZE variable in USOURCE(DASGUIDE).

³This design decision was made because there would be significant problems with interpreting the results from CPExpert if both types of analysis were performed simultaneously, and because a large coding effort would be required to perform the concurrent analysis. Please call Computer Management Sciences if you have a strong requirement for concurrent analysis.

Chapter 2.9: Exclude reporting low-activity data sets: **MIN42PCT** variable

If SMF Type 42 (Data Set Statistics) information is available⁴ in a MXG performance data base, the DASD Component will process MXG TYPE42DS to select data sets that were referenced during RMF measurement intervals in which the poorly performing devices exceeded the average performance. There can be many data sets referenced on the devices (hundreds or even thousands of data sets can be referenced). It is not helpful to have a large number of data sets listed. Additionally, some data sets might have a low I/O activity and would not be interesting to analyze. Consequently, CPExpert uses the **MIN42PCT** variable to limit the number of data sets listed in any RMF measurement interval, for a particular device. Only data sets having a percent of activity for the total volume greater than the percent specified by the MIN42PCT variable will be listed individually, and any remaining data sets will be summarized and listed on a single line.

The default value for the MIN42PCT variable is 0.1, indicating that data sets will not be listed individually for any RMF interval unless the data set intensity of access (I/O rate * response time) was greater than 0.1% of the total volume intensity of access. You can alter this percent of data sets listed by using the MIN42PCT guidance variable.

For example, if you specified **%LET MIN42PCT = 25;** in USOURCE(DASGUIDE), a maximum of 4 data sets would be listed in any RMF interval since no more than 4 data sets could have 25% or higher access intensity.

Please note that, regardless of the data set access intensity of any particular data set, only the number of data sets specified by the **LIST42DS** guidance variable (described earlier) will be listed. This means that there are two ways to limit the number of data sets listed: (1) the LIST42DS which limits the number of data sets listed in any RMF interval, and (2) the **MIN42PCT** guidance variable which limits the data sets listed to those that exceed the specified percent.

Chapter 2.10: List data for all RMF intervals: **LISTALL** variable

Once the DASD Component has selected one or more devices to analyze, it analyzes device-related data for all RMF intervals. This analysis determines whether each device has performance problems, and determines the most significant problems during the RMF interval. When the DASD Component produces the results of its analysis, it shows device information for each RMF interval analyzed, regardless of whether a performance problem existed with the device. This listing of all RMF intervals is the default, since it presents a continuous view of the device performance for all RMF intervals in the performance data base.

⁴%LET TYPE42DS = Y; was specified in USOURCE(GENGUIDE).

Some users analyze DASD problems for many systems in a single execution of the DASD Component, and they collect RMF data on short time intervals (for example, 15 minute RMF intervals or shorter). The default listing of data for all RMF intervals can produce very large reports for these users.

The LISTALL variable can be used to suppress the listing of RMF intervals in which no performance problem was detected with a device. You can suppress listing all RMF intervals by specifying **%LET LISTALL=N;** in USOURCE(DASGUIDE). With this specification, the DASD Component will list only those intervals in which performance problems were detected.

Chapter 2.11: Minimum I/O rate to analyze: MINIORT variable

Before the DASD Component selects a device for analysis, the code determines whether the device has an I/O rate greater than the MINIORT value. This feature allows user to exclude low-activity devices from analysis. The default value for the MINIORT variable is 0, indicating that all devices will be candidates for analysis, regardless of their I/O rate. You can alter this I/O rate by using the MINIORT guidance variable.

For example, if you specified **%LET MINIORT = 1;** in USOURCE(DASGUIDE), only devices with a minimum I/O rate greater than 1 I/O per second would be selected for analysis.

Chapter 2.12: Minimum I/O rate to analyze: MINIOWT variable

Before the DASD Component selects a device for analysis, the code determines whether the device has an I/O “intensity” (or weight) greater than the MINIOWT value. This feature allows user to exclude low-activity devices from analysis. The default value for the MINIOWT variable is 0, indicating that all devices will be candidates for analysis, regardless of their I/O intensity. You can alter this I/O rate by using the MINIOWT guidance variable.

For example, if you specified **%LET MINIOWT = 1000;** in USOURCE(DASGUIDE), only devices with a minimum I/O intensity greater than 1000 would be selected for analysis.

Chapter 2.13: Minimum I/O response to analyze: MINRESP variable

Before the DASD Component selects a device for analysis, the code determines whether the device has an I/O response greater than the MINRESP value. This feature allows user to exclude low-activity devices from analysis. The default value for the MINRESP variable is 0, indicating that all devices will be candidates for analysis, regardless of their I/O response. You can alter this I/O rate by using the MINRESP guidance variable. The MINRESP variable is specified in millisecond units.

For example, if you specified **%LET MINRESP = 5;** in **USOURCE(DASGUIDE)**, only devices with a minimum I/O response greater than milliseconds would be selected for analysis.

Chapter 2.14: Number of volumes to report: **REPORT** variable

The DASD Component analyzes performance information about the device(s) with the most performance improvement potential, based on **ANALYZE** guidance variable described earlier. To put the performance improvement potential of these devices in perspective relative to the overall sysplex or system, the DASD Component also produces a brief summary of relevant characteristics of other devices, with the list ranked in descending order by device access intensity (I/O RATE * average I/O response). The summary by sysplex is produced in Rule DAS000 and the summary by system is produced in Rule DAS050.

The number of devices for which summary information is produced is controlled by the **REPORT** guidance variable. The default for the **REPORT** guidance variable is 10, indicating that information is produced for the top 10 devices on a system basis. The value of the **REPORT** variable is doubled for the sysplex view.

Chapter 2.15: Selecting volumes to analyze: **SELECT** variable

The **SELECT** guidance variable is used to tell CPEXpert whether you wish to select specific volumes for analysis. The **SELECT** variable acts as a "switch" to control whether CPEXpert processes the **DASGUIDE** member searching for volumes to select for analysis. The point of having a "switch" variable is that some installations may wish to regularly analyze all volumes, but periodically select specific volumes for analysis.

If the **SELECT** variable is **N**, CPEXpert will not select specific volumes for analysis, but will analyze all volumes after applying data selection criteria. If the **SELECT** variable is **Y**, CPEXpert will process the **DASGUIDE** member to identify the volumes to be analyzed.

Please refer to Chapter 4 for more information on selecting volumes to analyze.

Chapter 2.16: Analyze shared DASD Conflicts - **SHARED** variable

CPEXpert can perform an analysis of conflicts between DASD shared between systems or MVS images. The shared DASD analysis is accomplished when you perform the following steps:

- You collect and process SMF data for each system to be analyzed. The SMF data must be placed into your standard performance data base. This step is not specific for

CPEXpert; it is the normal processing you would do. The step is listed first just for completeness.

- **You should ensure that the system clocks on all systems are set at roughly the same time.** CPEXpert will examine the SMF data contained in your performance data base, based upon the SMF time associated with the measurement data. Shared DASD problems between systems will be detected based on this SMF time. If the SMF times are significantly different, the analysis might not properly identify conflicts or conflicting applications.

You should not worry that the system clocks be set at **exactly** the same time. The analysis performed by CPEXpert is not intended to identify an isolated performance problem. Rather, CPEXpert attempts to identify those problems that **continually** cause shared DASD performance problems. If a shared DASD problem is continual, the problem will be reflected in multiple SMF recording intervals. Consequently, it is not essential that the SMF times be exactly matched between systems.

If the shared devices are in the same sysplex, this step will be performed automatically by the sysplex timer. The step is listed just for completeness.

- You specify **%LET SHARED = Y;** in USOURCE(DASGUIDE) to tell CPEXpert to perform analysis of shared DASD.
- If you wish the DASD Component to identify applications causing shared DASD conflicts, you must perform the following additional steps:
 - Install the modification to MXG or MICS to collect data set-related information as described in Section 2 of this User Manual.
 - Specify **%LET TYPE30DD = Y;** in the USOURCE(GENGUIDE) general guidance module. This tells CPEXpert that you have installed the modification to MXG or NeuMICS to collect application-related information.

Chapter 2.17: SMF Type 30 modification installed - TYPE30DD variable

The **TYPE30DD** guidance variable tells CPEXpert whether you have installed the CPEXpert modification to MXG or MICS code to collect SMF Type 30 (Data Definition) information⁵.

⁵The TYPE30DD guidance variable is normally specified when installing CPEXpert and is discussed in the *CPEXpert Installation Guide*. The TYPE30DD guidance variable is included in this document simply to remind you that the SMF Type 30 information is required if you wish to have CPEXpert show applications that access devices with poor performance, wish to perform "loved one" analysis, or wish to perform critical data set analysis. The TYPE30DD is not included in the USOURCE(DASGUIDE) member released with CPEXpert, since a specification in USOURCE(DASGUIDE) would override any specification in USOURCE(GENGUIDE).

This modification allows CPEXpert to relate DASD performance information contained in SMF Type 74 records to specific service classes or performance groups, and to relate the Type 74 information to specific jobs or job steps.

The TYPE30DD statement should be specified as **%LET TYPE30DD = Y**; if you have installed the CPEXpert modification to collect Type 30 (Data Definition) information for the system(s) being analyzed.

Chapter 2.18: SMF Type 42 (Data Set Statistics)- TYPE42DS variable

The **TYPE42DS** guidance variable tells CPEXpert whether SMF Type 42 (Data Set) records are available for analysis⁶. If you collect the SMF Type 42 (Data Set) records, CPEXpert can report data set information⁷ for data sets residing on volumes with poor performance.

The default value for the TYPE42DS guidance variable is “N”, indicating that the TYPE42DS data set is not available for analysis. You can change this guidance if Type 42 (Data Set) records are available, and if you wish CPEXpert to report data set information, by specifying **%LET TYPE42DS=Y**; in USOURCE(DASGUIDE).

⁶The TYPE42DS data set statistics information is available only with a MXG performance data base.

⁷The TYPE42DS guidance variable is normally specified when installing CPEXpert and is discussed in the *CPEXpert Installation Guide*. The TYPE42DS guidance variable is included in this document simply to remind you that the SMF Type 42 information is required if you wish to have CPEXpert show data sets that are accessed on devices with poor performance. The TYPE42DS is not included in the USOURCE(DASGUIDE) member released with CPEXpert, since a specification in USOURCE(DASGUIDE) would override any specification in USOURCE(GENGUIDE).

Chapter 3: Excluding volumes from analysis

For a variety of reasons, you may wish to exclude certain volumes from analysis. The most common reason is that CPExpert repeatedly identifies problems with particular volumes, but (1) you do not wish to make changes to correct the problems, (2) you are unable to make changes (because of application requirements or political realities), or (3) the DASD Component analysis is "flawed" because of data problems or data averaging. Whatever the reason, you may wish to exclude certain volumes from analysis.

Excluding volumes from analysis is optional; you do not have to exclude any volumes from analysis. If you do **not** exclude volumes from analysis, CPExpert will analyze every DASD VOLSER encountered in SMF Type 74 data after applying data selection criteria.

If you wish to exclude volumes from analysis, you must (1) use the **EXCLUDE** guidance variable and (2) define the volumes to exclude.

Chapter 3.1: EXCLUDE variable

The EXCLUDE guidance variable is used to tell CPExpert whether you wish to exclude volumes from analysis. The EXCLUDE variable acts as a "switch" to control whether CPExpert processes the DASGUIDE member searching for volumes to exclude from analysis. The point of having a "switch" variable is that some installations may wish to regularly exclude volumes from analysis, but periodically analyze all volumes.

If the EXCLUDE variable is **N**, CPExpert will not exclude any volumes from analysis. If the EXCLUDE variable is **Y**, CPExpert will process the DASGUIDE member to identify all volumes to be excluded.

Chapter 3.2: Defining volumes to exclude

You may exclude any number of volumes from analysis by coding the name of the volume in the "EXCLUDE VOLUME" portion of the DASGUIDE Module. You simply enter the VOLSER of each volume to be excluded after the "EXCLUDE" statement.

- The volumes to be excluded are entered one VOLSER per line or multiple VOLSERs per line.
- If you include multiple VOLSERs on a single line, you must separate them by blanks or commas.
- You can include comments on the line by placing "*" before the comment.

- You can use "generic" exclude logic by simply listing the first "n" characters of the VOLSERs you wish to exclude. For example, to exclude all VOLSERs beginning with ABC, simply specify ABC as the VOLSER to exclude.

CPEXpert will process the DASGUIDE Module, searching for the "EXCLUDE" statement. Any VOLSERs between the "EXCLUDE" statement and the "**/" statement will be placed into global SAS macro variables and may be excluded from analysis.

Exhibit 3-3 was extracted from the distributed USOURCE(DASGUIDE) module, and illustrates how to exclude volumes from analysis.

```
*****;
* OPTIONAL GUIDANCE TO EXCLUDE VOLUMES FROM ANALYSIS;
*****;
* DO NOT REMOVE OR ALTER THE FOLLOWING SAS MACRO COMMENT LINE!;
/* EXCLUDE THE FOLLOWING VOLSER FROM ANALYSIS:
EXCLUDE
MVS21AA    * SAMPLE VOLSER DEFINITION
CHKPT01    * SAMPLE VOLSER DEFINITION
CHPPT02    * SAMPLE VOLSER DEFINITION
SPOOL1,SPOOL2,SPOOL3
* DO NOT REMOVE THE FOLLOWING MACRO COMMENT LINE!
*/
```

EXCLUDING SPECIFIC VOLUMES FROM ANALYSIS

EXHIBIT 3-3

Chapter 4: Selecting specific volumes for analysis

You may wish to select specific volumes for analysis by the DASD Component, irrespective of whether these volumes are among the "worst" performing devices. For example, you may wish to examine only volumes containing specific application data sets.

Selecting volumes for analysis is optional; you do not have to select specific volumes for analysis. If you do not select specific volumes for analysis, CPExpert will analyze every DASD VOLSER encountered in SMF Type 74 data after applying data selection criteria. If you **do** select specific volumes for analysis, CPExpert will report on only those volumes (but CPExpert will analyze the performance of the volumes within the context of the overall I/O configuration).

If you wish to select volumes from analysis, you must (1) use the **SELECT** guidance variable and (2) define the volumes to select.

Chapter 4.1: SELECT variable

The SELECT guidance variable is used to tell CPExpert whether you wish to select specific volumes for analysis. The SELECT variable acts as a "switch" to control whether CPExpert processes the DASGUIDE member searching for volumes to select for analysis. The point of having a "switch" variable is that some installations may wish to regularly analyze all volumes, but periodically select specific volumes for analysis.

If the SELECT variable is **N**, CPExpert will not select specific volumes for analysis, but will analyze all volumes after applying data selection criteria. If the SELECT variable is **Y**, CPExpert will process the DASGUIDE member to identify the volumes to be analyzed.

Chapter 4.2: Defining volumes to analyze

You may select any number of volumes for specific analysis by coding the name of the volume in the "SELECT VOLSER" portion of the DASGUIDE Module. You simply enter the VOLSER of each volume to be excluded after the "SELECT" statement.

- The volumes to be selected are entered one VOLSER per line or multiple VOLSERs per line.
- If you include multiple VOLSERs on a single line, you must separate them by blanks or commas.
- You can include comments on the line by placing "*" before the comment.

- You can use "generic" select logic by simply listing the first "n" characters of the VOLSERs you wish to select. For example, to select all VOLSERs beginning with DEF, simply specify DEF as the VOLSER to select.

CPEXpert will process the DASGUIDE Module, searching for the "SELECT" statement. Any VOLSERs between the "SELECT" statement and the "**/" statement will be placed into global SAS macro variables and will be the only volumes CPEXpert analyzes in detail (note that CPEXpert will continue to process data describing the entire I/O configuration so that the analysis is done within the context of your configuration).

Exhibit 3-4 was extracted from the distributed USOURCE(DASGUIDE) module, and illustrates how to select specific volumes for analysis.

You should note that CPEXpert still applies the "worst volume" philosophy when it reports the results from the analysis. That is, only volumes exceeding the average of those selected will be analyzed in detail. You might be confused when you examine the "next worst volumes" report, unless you keep this in mind. (Give us feedback if you don't like this approach. If it is unsatisfactory, we can change the design.)

```
*****;
* OPTIONAL GUIDANCE TO SELECT SPECIFIC VOLUMES FOR ANALYSIS;
*****;
* DO NOT REMOVE OR ALTER THE FOLLOWING SAS MACRO COMMENT LINE!;
/* SELECT THE FOLLOWING VOLSER FOR ANALYSIS:
SELECT
PAGE21      * SAMPLE VOLSER DEFINITION
PAGE22      * SAMPLE VOLSER DEFINITION
PAGE23      * SAMPLE VOLSER DEFINITION
PAGE24      * SAMPLE VOLSER DEFINITION
* DO NOT REMOVE THE FOLLOWING MACRO COMMENT LINE!
*/
```

SELECTING SPECIFIC VOLUMES FOR ANALYSIS

EXHIBIT 3-4

Chapter 5: Analyzing response objectives for critical data sets

As described in Section 1 (Chapter 5.3: Expanded analysis - Specific Data Sets) and Section 2 (Chapter 4: Defining Critical Data Sets), you can identify critical data sets to the DASD Component and specify a response objective for these data sets. The DASD Component can then analyze the performance of your DASD configuration, in one of two ways: (1) analyzing data set response based on information from TYPE42DS statistics or (2) analyzing data based on information in TYPE14/15 and using the CPEXpert modification to MXG or NeuMICS as the SMF Type30 records are processed.

Chapter 5.1: Analysis based on TYPE42DS

This method can provide comprehensive information, without requiring the modification to MXG. This method **is not applicable** to NeuMICS, since NeuMICS does not retain sufficient information related to SMF Type 42 (Data Set Statistics).

MXG creates TYPE42DS from the SMF Type 42 (Subtype 6 - data set I/O statistics) records created by SMS. The TYPE42DS file contains I/O access characteristics information, at the data set level. CPEXpert extracts data set information and data set response statistics from TYPE42DS, and compares these to the response objectives you have made in USOURCE(DASGUIDE) about critical data sets. When any data set response time exceeds the specified objective, the DASD Component selects the data set and the volume it resides on for detailed analysis.

The following steps are necessary to implement the critical data set analysis based on TYPE42DS records:

- Specify **%LET DASDSN = Y;** in USOURCE(DASGUIDE) to tell the DASD Component to analyze DASD performance based upon data set name.
- Specify **%LET TYPE42DS = Y;** in the USOURCE(GENGUIDE) CPEXpert general guidance module. This tells CPEXpert that you have SMF Type 42 (Data Set Statistics) available in your performance data base.
- Define critical data sets to the DASD Component and specify a response objective for each data set. This process is described in Section 2 (Chapter 6).

Chapter 5.2: Analysis based on TYPE14/15 and CPEXpert modification

This method is a bit more involved, but can be used regardless of whether data sets are managed by SMS, and this method can be used regardless of whether you use MXG or NeuMICS to create your performance data base.

With this method, the DASD Component analyzes SMF Type 14/15 records to extract data set names that correspond to the critical data sets that you have identified in USOURCE(DASGUIDE). The CPExpert modification to MXG or NeuMICS extracts DD information as SMF Type 30 records are processed. The DASD Component then correlates the data set information with the DD information, to determine whether critical data sets exceed the specified response objective. When any data set responses time exceeds the specified objective, the DASD Component selects the data set and the volume it resides on for detailed analysis.

The following steps are necessary to implement the critical data set analysis based on SMF Type 14/15, Type 30, and Type 74

- Install the modification to MXG or NeuMICS to collect data set-related information as described in Section 2 of this User Manual. The DASD Component cannot perform expanded analysis based on specific data sets unless it has information relating DASD use to specific data sets, service classes, jobs, and job steps. This information is extracted from the SMF Type 30(DD) records by the modifications described in Section 2 (Chapter 3 and Chapter 4).
- Define critical data sets to the DASD Component and specified a response objective for each data set. This process is described in Section 2 (Chapter 6).
- Specify **%LET DASDSN = Y;** in USOURCE(DASGUIDE) to tell the DASD Component to analyze DASD performance based upon data set name.
- Specify **%LET TYPE30DD = Y;** in the USOURCE(GENGUIDE) CPExpert general guidance module. This tells CPExpert that you have installed the modification to MXG or MICS to collect application-related information.
- Collect SMF Type 14/15 records and execute the DAS1415 module (described in Section 4 of this manual) to process the SMF Type 14/15 records. You must execute the DAS1415 module **before** you execute your normal daily update of your performance data base.

Chapter 6: Analyzing VSAM data sets

VSAM data set activity typically accounts for a large percent of I/O activity (more than 70% at some sites). Tuning of a few files or correcting common problems often can result in significantly improved performance (IBM benchmarks show up to 90% improvement resulting from some simple changes).

Analysis of VSAM data sets is not applicable to performance data bases created with NeuMICS, since the required SMF Type 42(Data Set statistics) and Type 64 (VSAM statistics) are not available.

With CPExpert Release 12.2, the DASD Component was enhanced to provide a rudimentary analysis of common VSAM problems. Additional analysis will be added in future enhancements to the DASD Component.

The DASD Component can optionally analyze VSAM data set performance problems or potential problems, only if the MXG TYPE64 file is available (and CPExpert is provided with the **%LET TYPE64=Y**; guidance variable). Additionally, most analysis depends on having the MXG TYPE42DS file available (and CPExpert is provided with the **%LET TYPE42DS=Y**; guidance variable).

Most sites have many VSAM data sets. Some of the VSAM data sets are open for a long time, but most are open for only a short interval. Some of the VSAM data sets have a significant amount of I/O activity, while other VSAM data sets have very low activity. These characteristics can result in very large reports produced by the DASD Component, since the DASD component can report on problems with each VSAM data set.

Options are provided with the VSAM analysis to summarize findings for VSAM data sets, to ignore VSAM data sets that are open for a short time, or to ignore VSAM data sets that have less than a specified amount of I/O activity. Additionally, each finding is based on guidance variables. These guidance variables can be altered so CPExpert will produce results only when more restrictive conditions are met. All these options can be (and should be) used to reduce the output from the DASD Component's to a moderate size.

This chapter describes the guidance variables that are available for analyzing performance problems of VSAM data sets. Exhibit 3-5 illustrates the VSAM analysis guidance variables contained in USOURCE(DASGUIDE) Component.

```

*****;
*   VSAM ANALYSIS GUIDANCE VARIABLES                               ;
*****;
%LET ANALVSAM  =BAD;      * VSAM ANALYSIS OPTION                  ;
%LET CASPLITS  = 10;      * EXCESSIVE CONTROL AREA SPLITS        ;
%LET DIRINDEX  =25;      * PCT DIRECT ACCESS TO VSAM INDEX COMPONENT ;
%LET EXTENTS   = 1;      * EXCESSIVE SECONDARY EXTENTS            ;
%LET LSRSEQ    = 50;      * LSR SEQUENTIAL ACCESS DOMINATE        ;
%LET MXEXTENT  =225;      * MAXIMUM EXTENTS INDICATING POTENTIAL PROBLEM ;
%LET NSRDIR    = 50;      * NSR DIRECT ACCESS DOMINATE            ;
%LET OPENTIME  =900;      * MINIMUM VSAM DATA SET OPEN TIME FOR ANALYSIS ;
%LET PCTDIR    =80;      * PERCENT DIRECT VSAM ACCESSES          ;
%LET PCTSEQ    =80;      * PERCENT SEQUENTIAL VSAM ACCESSES      ;
%LET VSAMEXCL  = NO;      * OPTION: EXCLUDE SELECTED VSAM DATA SETS ;
%LET VSAMIO    = 100;     * VSAM I/O ACTIVITY SIGNIFICANT TO ANALYZE ;
%LET VSAMSMRY  = NO;      * VSAM DATA SET SUMMARIZED?            ;
*****;

```

VSAM ANALYSIS GUIDANCE VARIABLES

EXHIBIT 3-5

Chapter 6.1: Controlling analysis of VSAM: ANALVSAM variable

The DASD Component optionally analyzes VSAM data sets to identify performance problems or potential problems if SMF Type 64 information is available, and if **%LET TYPE64=Y**; has been specified in **USOURCE(DASGUIDE)**⁸.

For flexibility, there are four options that control whether and how the VSAM data set analysis is performed. These options are specified by the ANALVSAM guidance variable:

- **%LET ANALVSAM = BAD; is specified in USOURCE(DASGUIDE).** This is the default specification. When **%LET ANALVSAM=BAD;** is specified, the analysis is done **only** for VSAM data sets that reside on the devices that have been selected for analysis because the devices are the “worst performing” DASD devices in the configuration.
- **%LET ANALVSAM = ALL; is specified in USOURCE(DASGUIDE).** This option directs CPEXpert to analyze all VSAM data sets in the configuration, regardless of which

⁸ Most of the VSAM analysis also requires TYPE42DS information. Analysis that requires TYPE42DS information will be suppressed unless **%LET TYPE42DS=Y**; has been specified in either **USOURCE(GENGUIDE)** or **USOURCE(DASGUIDE)**.

devices the VSAM data sets reside. CPEXpert will analyze VSAM performance for all volumes, after analyzing those volumes whose performance is worse than average. The basic data selection criteria (for example, date and time, system, etc.) will be applied before performing the VSAM analysis.

- **%LET ANALVSAM = ONLY; is specified in USOURCE(DASGUIDE).** This option directs CPEXpert to analyze **only** VSAM data sets, based on SMF Type 64 and (optionally) SMF Type 42 (Data Set Statistics) records. When the ONLY option is selected, CPEXpert will analyze VSAM performance without performing the basic DASD analysis of worst performing volumes. Exercising this option eliminates the processing associated creating a model of the I/O configuration and analyzing device performance, and CPEXpert will ignore devices with non-VSAM data sets. This option will be particularly useful when new applications are produced and these applications use VSAM data sets.
- **%LET ANALVSAM = NO; is specified in USOURCE(DASGUIDE).** This option directs CPEXpert to eliminate the VSAM analysis. This option might be selected after CPEXpert has analyzed VSAM data set performance, and you have made all of the changes that you intend to implement.

Chapter 6.2: Excessive Control Area splits: CASPLITS variable

CPEXpert examines the SMF Type 64 information contained in the MXG TYPE64 data set to identify VSAM data sets that have excessive Control Area splits. CPEXpert sums the ACCASPLT variable (the number of CA splits since the data set was created) and the CASPLITS variable (the number of CA splits with the current OPEN). CPEXpert compares this sum with the **CASPLITS** guidance variable in USOURCE(DASGUIDE). CPEXpert produces Rule DAS600 when the total number of CA splits exceeds the value specified by the **CASPLITS** guidance variable.

The default value for the **CASPLITS** guidance variable is 10, indicating that CPEXpert should produce Rule DAS600 when a VSAM data set experienced more than 10 CA splits. You should alter this guidance variable if you feel that Rule DAS600 is produced too often, or if you do not wish to take action when only 10 CA splits occur. For example, if you wish CPEXpert to produce Rule DAS600 only when more than 50 CA splits occur for any VSAM data set, specify **%LET CASPLITS=50;** in USOURCE(DASGUIDE).

Chapter 6.3: Percent direct access to VSAM index component: DIRINDEX variable

CPEXpert examines the SMF Type 64 information contained in MXG TYPE64 data set to identify VSAM KSDS or VRRDS data sets that have insufficient buffers assigned to the index component. CPEXpert uses the TYPE42DS information to compute the percent of direct accesses to the VSAM data set, using the following algorithm:

$$\text{Percent direct accesses} = \frac{S42AMDRB}{S42AMSRB \% S42AMDRB}$$

where: S42AMSRB = Blocks read using sequential access

S42AMDRB = Blocks read using direct access

CPEXpert produces *Rule DAS621* under the following conditions:

- The TYPE42DS S42DSBUF variable showed that NSR was used for KSDS or VRRDS VSAM data sets, **and**
- The percent of direct accesses for the index component was greater than **DIRINDEX** guidance variable in USOURCE(DASGUIDE), **and**
- The number of buffers (the MXG BUFDRNO variable) assigned to the index component was less than the number of index levels (the MXG ACCLEVEL variable) for the VSAM data set.

CPEXpert produces *Rule DAS622* under the following conditions:

- The TYPE42DS S42DSBUF variable showed that NSR was used for KSDS or VRRDS VSAM data sets, **and**
- The percent of direct accesses for the index component was greater than **DIRINDEX** guidance variable in USOURCE(DASGUIDE), **and**
- The STRNO specification in the ACB (the MXG ACBSTRNO) was greater than one (indicating that concurrent accesses had been specified for direct processing), **and**
- The number of buffers (the MXG BUFDRNO variable) assigned to the index component was less than the ACBSTRNO value, plus 1.

The default value for the DIRINDEX guidance variable is 25%, so CPEXpert will produce Rule DAS621 or DAS622 for NSR VSAM data sets when more than 25% of the accesses were direct for the index component and the number of buffers assigned to the index component was less than the number of index levels (for rule DAS621) or less than the

STRNO value in the ACB, plus 1 (for Rule DAS622). You can change the DIRINDEX guidance variable if you feel that Rule DAS621 or Rule DAS622 are produced too often, or if you do not wish to take action when to increase the number of buffers assigned to the VSAM data sets listed.

For example, if you wish CPEXpert to produce Rule DAS621 or Rule DAS622 only when more than 50% of the accesses for an index component were direct, specify **%LET DIRINDEX=50;** in USOURCE(DASGUIDE). You can completely suppress Rule DAS621 and Rule 622 by specifying **%LET DIRINDEX=100;** in USOURCE(DASGUIDE), since the percent cannot exceed 100.

Chapter 6.4: Excessive EXTENTS were allocated: EXTENTS variable

CPEXpert examines the SMF Type 64 information contained in MXG TYPE64 data set to identify VSAM KSDS or VRRDS data sets that have excessive secondary allocations.

CPEXpert compares NREXTNTS variable (the number of secondary extents in the VSAM data set this OPEN) with the **EXTENTS** guidance variable in USOURCE(DASGUIDE). CPEXpert produces Rule DAS604 when the NREXTENT (the total number of extents) is greater than one, and the number of secondary extents allocated for this OPEN exceeds the value specified by the EXTENTS guidance variable.

The default value for the EXTENTS guidance variable is 0, indicating that CPEXpert should produce Rule DAS604, Rule DAS605, or Rule DAS606 (depending on data encountered) when any secondary extent was allocated for the VSAM data sets listed. You can change the EXTENTS guidance variable if you feel that Rule DAS604 is produced too often, or if you do not wish to take action when secondary extents are allocated. For example, if you wish CPEXpert to produce Rule DAS604 only when more than 5 secondary extents were allocated for VSAM data sets, specify **%LET EXTENTS=5;** in USOURCE(DASGUIDE). You can completely suppress Rule DAS635 by specifying **%LET EXTENTS=255;** in USOURCE(DASGUIDE), since VSAM cannot allocate 255 extents.

Chapter 6.5: Specifying LSR sequential access domination: LSRSEQ variable

With Local Shared Resource (LSR), VSAM buffers normally are shared among VSAM data sets accessed by tasks in the same address space. Since LSR is oriented toward shared (and direct) access, there is an expectation that a record might be re-used. Consequently, buffer management algorithms retain buffers as long as possible, using a least-recently used (LRU) algorithm, after a record is processed from the LSR buffers. LSR is not suited for applications that use sequential or skip sequential as their primary access mode, because there is no read-ahead algorithm with LSR, and there is no inherent overlap of I/O

and CPU processing. Consequently, using LSR for sequential access processing could degrade rather than improve performance.

After applying the screening criteria specified for VSAM data sets, and extracting SMF Type 64 information for those VSAM data sets, CPEXpert examines SMF Type 42 (Data Set Statistics) information for the selected VSAM data sets. CPEXpert uses the TYPE42DS information to compute the percent of sequential accesses to the VSAM data set, using the following algorithm:

$$\text{Percent sequential accesses} = \frac{S42AMSRB}{S42AMSRB \% S42AMDRB}$$

where: S42AMSRB = Blocks read using sequential access

S42AMDRB = Blocks read using direct access

CPEXpert produces Rule DAS635 under the following conditions:

- The TYPE42DS S42DSBUF variable showed that LSR was used for KSDS or VRRDS VSAM data sets, and
- The percent of sequential accesses for the data component was greater than the value specified for the **LSRSEQ** guidance variable in USOURCE(DASGUIDE).

The default value for the LSRSEQ guidance variable is 75%, so CPEXpert will produce Rule DAS635 when LSR was specified as the buffering technique, and more than 75% of the accesses were sequential for the data component.

You can change the percent of sequential access that CPEXpert uses to determine whether to produce Rule DAS635 by altering the LSRSEQ guidance variable. For example, if you wish CPEXpert to produce Rule DAS635 only when more than 90% of the accesses to a LSR buffer pool were sequential, specify **%LET LSRSEQ=90;** in USOURCE(DASGUIDE). You can completely suppress Rule DAS635 by specifying **%LET LSRSEQ=100;** in USOURCE(DASGUIDE), since the percent cannot exceed 100.

Chapter 6.6: Specifying maximum extents: MXEXTENT variable

CPEXpert examines the SMF Type 64 information contained in the MXG TYPE64 data set to identify VSAM data sets are in danger of reaching the maximum allowed number of extents. CPEXpert compares NREXTENT variable (the total number of extents in the VSAM data set) with the **MXEXTENT** guidance variable in USOURCE(DASGUIDE). CPEXpert produces Rule DAS607 when the NREXTENT is greater than the value specified by the MXEXTENT guidance variable **and** at least one extent was allocated during the

current OPEN of the data set (CPEXpert uses the NREXTNTS variable in TYPE64 for this decision).

The default value of the MXEXTENT guidance variable is 225, indicating that CPEXpert should produce Rule DAS607 when at least 225 extents have been allocated for a VSAM data set. Since the maximum allowable is 255, the default value provides a threshold at which CPEXpert provides notification that there is a potential problem.

You can alter this threshold by using the MXEXTENT guidance variable. For example, if you wish to be notified when the number of extents reach 200 (and at least one extent was required with the current OPEN), specify **%LET MXEXTENT=200;** in USOURCE(DASGUIDE).

Chapter 6.7: Specifying NSR direct access domination: NSRDIR variable

Non-shared resource (NSR) is the default VSAM buffering technique. VSAM data sets with NSR buffering can be accessed sequentially or direct (or both). However, NSR is suited for sequential processing because, if the data set access is sequential, the buffers are managed with a read-ahead algorithm. The read-ahead algorithm provides overlap of I/O and CPU processing and is efficient for sequential accesses. Since NSR is oriented toward sequential access, there is no expectation that a record will be re-used (as might exist with direct processing). Consequently, once a record is processed from the NSR buffers, the buffer is likely to be reclaimed for another record read from DASD.

NSR is not suited for direct processing, although NSR often is used for direct processing because it is easy to use and is the default buffering technique. Nonetheless, performance can be significantly improved if LSR is used for direct processing of VSAM data sets.

After applying the screening criteria specified for VSAM data sets, and extracting SMF Type 64 information for those VSAM data sets, CPEXpert examines SMF Type 42 (Data Set Statistics) information for the selected VSAM data sets. CPEXpert uses the TYPE42DS information to compute the percent of direct accesses to the VSAM data, using the following algorithm:

$$\text{Percent direct accesses} = \frac{S42AMDRB}{S42AMSRB \% S42AMDRB}$$

where: S42AMSRB = Blocks read using sequential access

S42AMDRB = Blocks read using direct access

CPEXpert produces Rule DAS625 under the following conditions:

- The TYPE42DS S42DSBUF variable showed that NSR was used for KSDS or VRRDS VSAM data sets, and
- The percent of direct accesses for the data component was greater than the value specified for the NSRDIR guidance variable in USOURCE(DASGUIDE).

The default value for the NSRDIR guidance variable is 75%, so CPExpert will produce Rule DAS625 when NSR was specified as the buffering technique, and more than 75% of the accesses were direct for the data component.

You can change the percent of sequential access that CPExpert uses to determine whether to produce Rule DAS635 by altering the NSRDIR guidance variable. For example, if you wish CPExpert to produce Rule DAS625 only when more than 90% of the accesses to a VSAM data set with NSR were sequential, specify **%LET NSRDIR=90;** in USOURCE(DASGUIDE). You can completely suppress Rule DAS625 by specifying **%LET NSRDIR=100;** in USOURCE(DASGUIDE), since the percent cannot exceed 100.

Chapter 6.8: Minimum VSAM open time: OPENTIME variable

Many installations have a large number of VSAM data sets that are open for a relatively short time, and analyzing these data sets would unnecessarily clutter the reports produced by the DASD Component. The **OPENTIME** guidance variable allows installations to control which data sets are analyzed, based on the amount of elapsed time that the data sets are open.

The SMF variables reflecting when the data set was opened (SMF30TM and SMF30DT) were available with z/OS Version 1 Release 1 (z/OS V1R1). For versions of MVS prior to z/OS V1R1, the default value for the OPENTIME guidance variable is zero, and this value cannot be changed (since CPExpert would have nothing in SMF to compare against). Beginning with z/OS V1R1, the default value for the OPENTIME guidance variable is 300, indicating that CPExpert should ignore⁹ VSAM data sets that are open for less than 300 seconds (or 5 minutes).

If you are running MVS with z/OS V1R2 or subsequent release, and you wish CPExpert to ignore VSAM data sets that are open for different elapsed time, you can change the OPENTIME guidance variable to a different number of seconds. For example, if you wish CPExpert to analyze only those VSAM data sets that are open for more than 30 minutes, specify **%LET OPENTIME=1800;** in USOURCE(DASGUIDE).

⁹CPExpert does NOT ignore any VSAM data set based on open time if **%LET VSAMSMRY=Y;** has been specified in USOURCE(DASGUIDE), until all activity for each VSAM data set has been summarized. If **%LET VSAMSMRY=Y;** has been specified, the guidance is applied to the summarized data.

Chapter 6.9: Specifying percent direct access for control interval size: PCTDIR variable

A *Control Interval* is a continuous area of direct access storage that VSAM uses to store logical records. The size of Control Intervals can vary from one VSAM data set to another, but all the Control Intervals within the data portion of a particular data set must be the same length. The type of processing that is used should guide the choice of control interval size. When direct processing accounts for most of the accesses, a small data Control Interval is preferable. This is because only one record is retrieved at a time with direct processing. If a large Control Interval is specified, unnecessary I/O overhead is incurred reading the excess information. IBM suggests that a 4096 byte Control Interval normally would be appropriate for direct processing.

After applying the screening criteria specified for VSAM data sets, and extracting SMF Type 64 information for those VSAM data sets, CPEXpert examines SMF Type 42 (Data Set Statistics) information for the selected VSAM data sets. CPEXpert uses the TYPE42DS information to compute the percent of accesses to the data component of the VSAM data set that were direct, using the following algorithm:

$$\text{Percent sequential accesses} = \frac{S42AMSRB}{S42AMSRB \% S42AMDRB}$$

where: S42AMSRB = Blocks read using sequential access

S42AMDRB = Blocks read using direct access

CPEXpert produces Rule DAS611 when the percent of direct accesses for the data component was greater than the **PCTDIR** guidance variable in USOURCE(DASGUIDE), and the Control Interval size was greater than 4096 bytes. Additionally, CPEXpert verifies that the maximum logical records (maximum LRECL) is less than 50% of the Control Interval size. This verification is done to make sure that Rule DAS611 is not produced for VSAM data sets that have spanned records.

The default value for the PCTDIR guidance variable is 80%, so CPEXpert will produce Rule DAS611 when more than 80% of the accesses were direct for the data component, and the Control Interval size was more than 4096 bytes. You can alter this algorithm by changing the PCTDIR guidance variable in USOURCE(DASGUIDE). For example, if you wish CPEXpert to produce Rule DAS611 only when more than 90% of the accesses to the data component were direct, specify **%LET PCTDIR=90;** in USOURCE(DASGUIDE). You can completely suppress Rule DAS611 by specifying **%LET PCTDIR=100;** in USOURCE(DASGUIDE), since the percent cannot exceed 100.

Chapter 6.10: Specifying percent sequential access for control interval size: PCTSEQ variable

A *Control Interval* is a continuous area of direct access storage that VSAM uses to store logical records. The size of Control Intervals can vary from one VSAM data set to another, but all the Control Intervals within the data portion of a particular data set must be the same length. The type of processing that is used should guide the choice of control interval size. When sequential processing accounts for most of the accesses, a large data Control Interval would normally be a good choice. This is because multiple records can be read into buffers and processed sequentially. For example, given a 16KB data buffer space, it is better to read two 8 KB Control Intervals with one I/O operation than four 4 KB Control Intervals with two I/O operations.

After applying the screening criteria specified for VSAM data sets, and extracting SMF Type 64 information for those VSAM data sets, CPEXpert examines SMF Type 42 (Data Set Statistics) information for the selected VSAM data sets. CPEXpert uses the TYPE42DS information to compute the percent of sequential accesses to the data component of the VSAM data set, using the following algorithm:

$$\text{Percent sequential accesses} = \frac{S42AMSRB}{S42AMSRB \% S42AMDRB}$$

where: S42AMSRB = Blocks read using sequential access

S42AMDRB = Blocks read using direct access

CPEXpert produces Rule DAS610 when the percent of sequential accesses for the data component was greater than the **PCTSEQ** guidance variable in USOURCE(DASGUIDE), and the Control Interval size was less than 8192 bytes.

The default value for the PCTSEQ guidance variable is 80%, so CPEXpert will produce Rule DAS610 when more than 80% of the accesses were sequential for the data component, and the Control Interval size was less than 8192 bytes. You can alter this algorithm by changing the PCTSEQ guidance variable in USOURCE(DASGUIDE). For example, if you wish CPEXpert to produce Rule DAS610 only when more than 90% of the accesses to the data component were sequential, specify **%LET PCTSEQ=90;** in USOURCE(DASGUIDE). You can completely suppress Rule DAS610 by specifying **%LET PCTSEQ=100;** in USOURCE(DASGUIDE), since the percent cannot exceed 100.

Chapter 6.11: Excluding VSAM data sets: VSAMEXCL variable

The **VSAMEXCL** guidance variable tells CPEXpert whether you wish to exclude certain VSAM data sets from analysis. The VSAMEXCL variable acts as a "switch" to control whether CPEXpert processes the DASGUIDE member searching for VSAM data

set names to exclude from analysis. The point of having a "switch" variable is that some installations may wish to regularly exclude VSAM data sets from analysis, but periodically analyze all VSAM data sets.

If the VSAMEXCL variable is **N**, CPEXpert will not exclude any VSAM data sets from analysis¹⁰. If the VSAMEXCL variable is **Y**, CPEXpert will process the DASGUIDE member to identify all VSAM data sets to be excluded.

Please refer to Chapter 7 for more information on excluding VSAM data sets from analysis.

Chapter 6.12: Specifying significant VSAM I/O activity: VSAMIO variable

Many installations have a large number of VSAM data sets that have little I/O activity, and analyzing these data sets would unnecessarily clutter the reports produced by the DASD Component. The **VSAMIO** guidance variable allows installations to control which data sets are analyzed, based on the amount of I/O activity of the data sets. CPEXpert examines the EXCP variable (number of EXCPs) in MXG TYPE64 records. CPEXpert ignores¹¹ all records that have less than the amount of I/O activity specified in the VSAMIO guidance variable.

The default value for the VSAMIO guidance variable is 100, indicating that CPEXpert should ignore VSAM data sets that have less than 100 I/O operations. This low value was selected as a default so that you could appreciate potential problems with most VSAM data sets in your installation.

If you wish CPEXpert to ignore VSAM data sets that have a different I/O activity, you can change the VSAMIO guidance variable. For example, if you wish CPEXpert to analyze only those VSAM data sets that have more than 1000 EXCPs during each OPEN, specify **%LET VSAMIO=1000;** in USOURCE(DASGUIDE).

¹⁰ Please note that the ANALVSAM guidance variable also acts to exclude VSAM data sets from analysis. If **%LET ANALVSAM=BAD;** is specified in USOURCE(DASGUIDE), only VSAM data sets residing on poorly performing volumes will be analyzed.

¹¹ CPEXpert does NOT ignore any VSAM data set based on I/O activity if **%LET VSAMSMRY=Y;** has been specified in USOURCE(DASGUIDE), until all activity for each VSAM data set has been summarized. If **%LET VSAMSMRY=Y;** has been specified, the guidance is applied to the summarized data.

Chapter 6.13: Summarizing VSAM activity: VSAMSMRY variable

Many installations have a large number of VSAM data sets that are open for a very short time, or have VSAM data sets that exhibit little I/O activity. Analyzing these data sets would unnecessarily clutter the reports produced by the DASD Component and produce a large amount of output¹². The **VSAMSMRY** guidance variable can be used to cause CPEXpert to summarize VSAM findings.

If **%LET VSAMSMRY=N;** is specified in USOURCE(DASGUIDE), CPEXpert will list every occasion in which a finding applies to a specific VSAM data set. This listing will show the time that the VSAM information was written to SMF, based on merging of MXG TYPE64 and TYPE42DS records. The information will show the JOB NAME using the VSAM data set, and other information associated with the rule. Since many VSAM data sets are OPENed and CLOSEd frequently (particularly with CLOSE=T), a very large number of lines of output can be produced. However, this level of output might be useful if VSAM problems are to be correlated with specific times of poor performance.

If **%LET VSAMSMRY=Y;** is specified in USOURCE(DASGUIDE), CPEXpert will summarize information for each VSAM data set and apply the VSAM analysis algorithms to the summarized information. The listing will show the time that the **last** VSAM information was written to SMF, based on merging of MXG TYPE64 and TYPE42DS records. The information will show the **last** JOB NAME using the VSAM data set, and a summary of other information associated with the rule. This level of output is useful for analyzing whether problems occur with VSAM data sets, over the entire measurement interval retained in the performance data base.

¹²Hundreds of pages of output might be created by CPEXpert if every VSAM data set were analyzed each time that it was OPENed.

Chapter 7: Excluding VSAM data sets from analysis

For a variety of reasons, you may wish to exclude certain VSAM data sets from analysis. The most common reason is that CPExpert repeatedly identifies problems with particular VSAM data sets, but (1) you do not wish to make changes to correct the problems, (2) you are unable to make changes (because of application requirements or political realities), or (3) the DASD Component analysis is "flawed" because of data problems or data averaging. Whatever the reason, you may wish to exclude certain VSAM data sets from analysis.

Excluding VSAM data sets from analysis is optional; you do not have to exclude any VSAM data sets from analysis. If you do **not** exclude VSAM data sets from analysis, CPExpert will analyze every VSAM data set encountered in SMF Type 64 after applying data selection criteria.

If you wish to exclude VSAM data sets from analysis, you must (1) use the **VSAMEXCL** guidance variable and (2) define the VSAM data sets to exclude.

Chapter 7.1: VSAMEXCL variable

The VSAMEXCL guidance variable is used to tell CPExpert whether you wish to exclude VSAM data sets from analysis. The VSAMEXCL variable acts as a "switch" to control whether CPExpert processes the DASGUIDE member searching for VSAM data sets to exclude from analysis. The point of having a "switch" variable is that some installations may wish to regularly exclude VSAM data sets from analysis, but periodically analyze all VSAM data sets.

If the VSAMEXCL variable is **N**, CPExpert will not exclude any VSAM data sets from analysis. If the VSAMEXCL variable is **Y**, CPExpert will process the DASGUIDE member to identify all VSAM data sets to be excluded.

Chapter 7.2: Defining VSAM data sets to exclude

You may exclude any number of VSAM data sets from analysis by coding the name of the VSAM data set in the "EXCLUDE SPECIFIC VSAM DATA SETS FROM ANALYSIS" portion of the DASGUIDE Module. You simply enter the VSAM data set name of each VSAM data set to be excluded after the `/* EXCLUDE VSAM DATA SETS FROM ANALYSIS` statement.

- The VSAM data sets to be excluded are entered one data set name per line or multiple data set names per line.
- If you include multiple data set names on a single line, you must separate them by blanks or commas.

- You can include comments on the line by placing " *" before the comment. NOTE that a blank character MUST precede the asterisk.
- You can use "generic" exclude logic by simply listing the first "n" characters of the VSAM data set names you wish to exclude, followed by an asterisk. For example, to exclude all VSAM data set names beginning with "D10.RMF.MONITOR3" simply specify "D10.RMF.MONITOR3*" as the data set name to exclude.

CPEXpert will process the DASGUIDE Module, searching for the "/* EXCLUDE VSAM DATA SETS FROM ANALYSIS" statement. Any VSAM data sets between the "/* EXCLUDE VSAM DATA SETS FROM ANALYSIS" statement and the "*/" statement will be placed into global SAS macro variables and may be excluded from analysis.

Exhibit 3-6 was extracted from the distributed USOURCE(DASGUIDE) module, and illustrates how to exclude VSAM data sets from analysis.

```

*****;
* OPTIONAL GUIDANCE TO EXCLUDE SPECIFIC VSAM DATA SETS FROM ANALYSIS ;
*****;
* DO NOT REMOVE OR ALTER THE FOLLOWING SAS MACRO COMMENT LINE! ;
/* EXCLUDE VSAM DATA SETS FROM ANALYSIS
DSN=RLSADSW.VF05D.ITEMACT.DATA
DSN=RLSADSW.VF07D.ITEMACT
DSN=D10.RMF.MONITOR3*
* DO NOT REMOVE THE FOLLOWING MACRO COMMENT LINE!
*/
*****;

```

EXCLUDING VSAM DATA SETS FROM ANALYSIS

EXHIBIT 3-6

Section 4: Executing the DASD Component

This section describes how to execute the DASD Component of CPEXpert.

The instructions in this section assume that you have installed the CPEXpert software. The DASD Component is installed as normal part of installing CPEXpert, and the modification to MXG or to MICS is installed as described in Section 2 of this User Manual. If you have not installed CPEXpert and the modification to MXG or NeuMICS, please install the software before continuing.

Executing the DASD Component involves executing the DASCPE Module. Additionally, the DAS1415 Module must be executed if you wish to analyze the DASD performance provided to critical data sets.

Chapter 1: Executing the DASCPE Module

This chapter describes how to execute the DASCPE Module of the DASD Component.

As stated in the Introduction to this document, the DASD Component consists of numerous modules working together to (1) shape system performance and utilization data for detailed analysis by other modules, (2) evaluate the data to assess potential causes of performance, (3) describe the results from the evaluation, and (4) maintain a historical record of the results from the analysis. These modules are loaded and controlled by the central DASD Component of CPEXpert (titled DASCPE).

Step 1. Use TSO ISPF to change the "prefix" in the data set names

Use TSO ISPF to change the "prefix" in the data set names (DSN) in the **USOURCE** DD statement, the **SOURCE** DD statement, the **CPEDATA** DD statement, the **CPEDASD** DD statement, the **HISTORY** DD statement, the **PDBLIB** DD statement, and the **SYSIN** DD statement of the JCL in accordance with your installation standards. The JCL is illustrated in Exhibit 4-1. (A "shell" of this JCL is contained on the distribution tape as "DASJCL2")

//jobname	JOB	job card information
//STEP02	EXEC	SAS,OPTIONS='MACRO DQUOTE PAGESIZE=65 ERRORABEND'
//USOURCE	DD	DSN=prefix.CPEXPRT.USOURCE,DISP=SHR
//SOURCE	DD	DSN=prefix.CPEXPRT.SOURCE,DISP=SHR
//CPEDATA	DD	DSN=prefix.CPEXPRT.CPEDATA,DISP=OLD
//CPEDASD	DD	DSN=prefix.CPEXPRT.CPEDATA,DISP=OLD
//PDBLIB	DD	DSN=prefix.MXG.MON,DISP=SHR
//LIBRARY	DD	DSN=saslib containing MXG FORMATS
//SYSIN	DD	DSN=prefix.CPEXPRT.SOURCE(DASCPE),DISP=SHR

JOB CONTROL LANGUAGE TO EXECUTE THE DASCPE MODULE

EXHIBIT 4-1

The CPEDATA DD statement in Exhibit 4-1 refers to the SAS data library maintained by CPExpert. The space for this library was created during the installation of CPExpert.

The CPEDASD DD statement in Exhibit 4-1 refers to the SAS library containing the DASD information created by the CPExpert modification to MXG or MICS. The CPEDASD DD statement is required only if you have installed the modification to MXG or MICS necessary to allow CPExpert to collect DASD I/O information at the job step level.

The PDBLIB DD statement in Exhibit 4-1 refers to the SAS library containing the performance data base to be analyzed. The example shows a sample DSN for a typical MXG performance data base. The DSN would be changed to "DSN=prefix.RMF.MICS.DETAIL" to use a MICS performance data base.

Exhibit 4-1 does not show the optional DD statements for MICS Information Areas (i.e., BATLIB DD, SCBLIB DD, HARLIB DD, etc.). The *CPExpert Installation Guide* describes how to use these optional DD statements if you have your MICS performance data base separated by MICS Information Area.

The LIBRARY DD statement in Exhibit 4-1 refers to the SAS library containing the MXG FORMATS. If you execute CPExpert against a MICS performance data base, the SASLIB DD statement would be changed to refer to the MICS FORMAT library. (CPExpert actually does not use the MICS FORMATS. You can specify the SAS **NOFMterr** option to eliminate problems with SAS FORMAT errors.)

Step 2: Make any appropriate changes to the DASGUIDE Module

Before submitting the JCL shown in Exhibit 4-1 and executing the DASDCPE Module, you should make appropriate changes to the CPEXPERT.USOURCE(DASGUIDE) module. These changes are described in Section 3 of this manual.

Step 3. Execute the DASCPE Module

Submit the JCL to execute the DASCPE Module. Most installations execute the DASCPE Module on a daily basis, after their normal update of their performance data base.

Chapter 2: Executing the DAS1415 Module

This section describes how to execute the DAS1415 Module of the DASD Component. The DAS1415 Module is executed **only** if you wish to perform expanded analysis based on data set name.

The DAS1415 Module reads the SMF Type 14/15 records to acquire information necessary to identify the DD names associated with data set names defined in USOURCE(DASGUIDE). As explained in Section 1 (Chapter 5.3), the DD names are later used in a modification to MXG or MICS processing to identify the activity of devices associated with the DD names.

The DAS1415 Module must be executed before MXG or MICS are executed to perform their normal daily update of your performance data base. This is necessary because the modification to MXG or MICS reads the data set created by the DAS1415 Module. [The data set created by the DAS1415 Module is simply the output from a SAS PROC FREQ creating a distribution of unique DD names which were used to reference data sets defined in USOURCE(DASGUIDE).]

Step 1. Use TSO ISPF to change the DD statements

Use TSO ISPF to change the SMF DD statement to refer to the current SMF data. Use TSO ISPF to change the "prefix" in the USOURCE DD statement, the SOURCE DD statement, the CPEDASD DD statement, and the SYSIN DD statement of the JCL in accordance with your installation standards. The JCL is illustrated in Exhibit 4-2. (A "shell" of this JCL is contained on the distribution tape as "DASJCL3")

//jobname	JOB	job card information
//STEP1	EXEC	SAS,OPTIONS='MACRO DQUOTE PAGESIZE=65 ERRORABEND'
//SMF	DD	DSN=data set name of current SMF file,DISP=SHR
//USOURCE	DD	DSN=prefix.CPEXPRT.USOURCE,DISP=OLD
//SOURCE	DD	DSN=prefix.CPEXPRT.SOURCE,DISP=SHR
//LIBRARY	DD	DSN=saslib containing MXG FORMATs
//CPEDASD	DD	DSN=prefix.CPEXPRT.CPEDASD,DISP=OLD
//SYSIN	DD	DSN=prefix.CPEXPRT.SOURCE(DAS1415),DISP=SHR

JOB CONTROL LANGUAGE TO EXECUTE THE DAS1415 MODULE

EXHIBIT 4-2

Step 2. Execute the DAS1415 Module

The next step is to execute the DAS1415 Module to process the USOURCE(DASGUIDE) information and to process the SMF Type 14/15 data.

The DAS1415 Module will create two SAS data sets in the CPEDASD SAS library, containing information extracted from USOURCE(DASGUIDE) and containing information extracted from the SMF Type 14/15 records for data sets matching the data set names defined in USOURCE(DASGUIDE).

Checklist for Executing the DASD Component, Mainframe

- ☐ Execute the DASD Component.
 - ☐ Change the "prefix" in the data set names in the DD statements.
 - ☐ Make any necessary changes to the DASGUIDE Module in USOURCE.
 - ☐ Submit the JCL to execute the DASDCPE Module.

Checklist for Executing DASD Component, Personal Computer

- ☐ Execute the DASD Component.
 - ☐ Identify the directories to SAS and CPExpert.
 - ☐ USOURCE filename at SAS PGM window
 - ☐ SOURCE filename at SAS PGM window
 - ☐ LIBRARY filename at SAS PGM window
 - ☐ Make any necessary changes to the DASGUIDE Module in USOURCE.
 - ☐ If you are executing under Windows, enter "%INCLUDE SOURCE(DASCPE)" at SAS PGM window and submit.
 - ☐ If you are executing under OS/2, enter "%INCLUDE SOURCE(DASCPE.SAS)" at SAS PGM window and submit.

Checklist for Performing Expanded Analysis

This checklist contains additional steps which must be performed if you wish the DASD Component to perform expanded analysis.

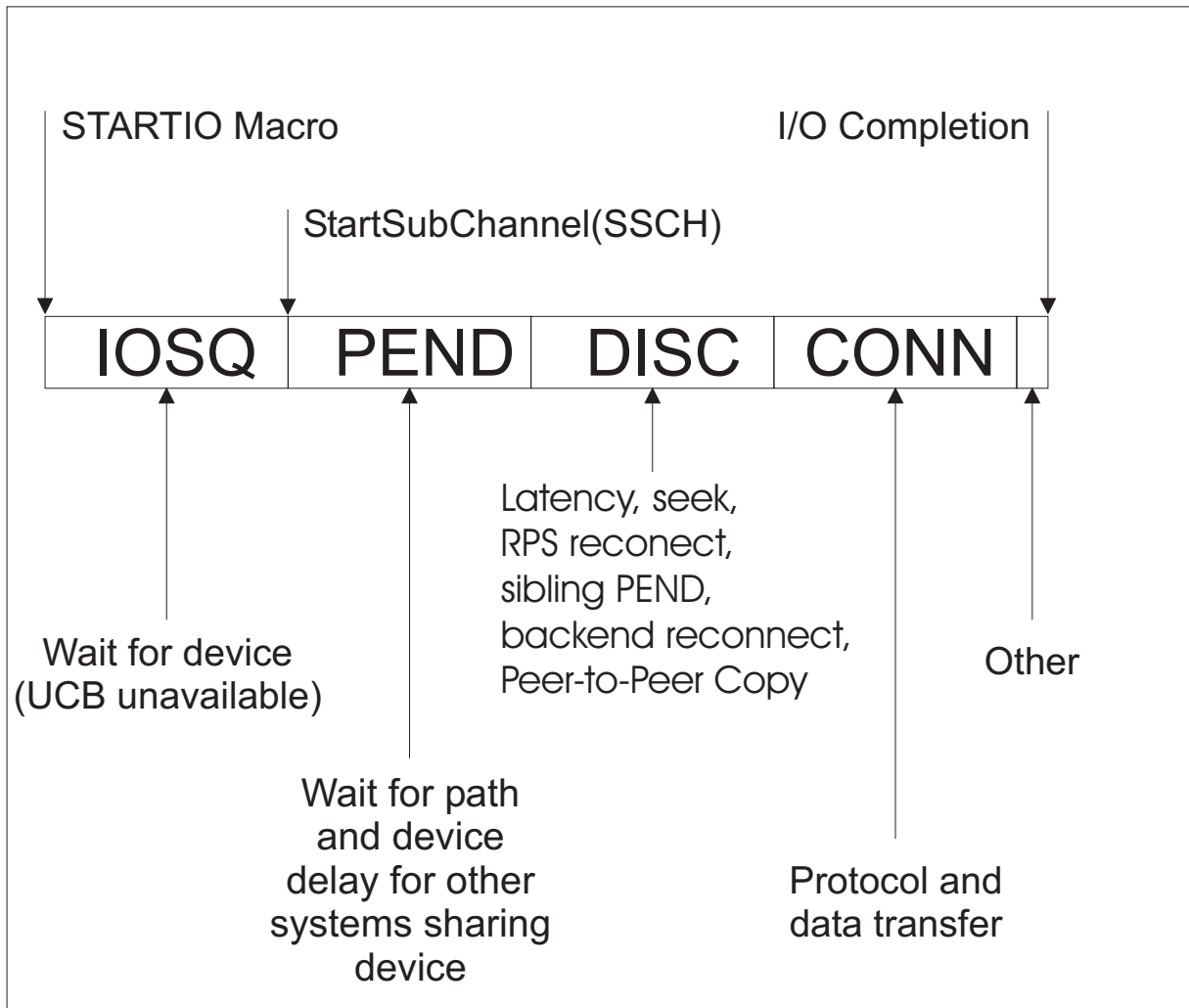
- ☐ Make sure that the CPExpert modification to MXG or to MICS is properly installed.
- ☐ Specify **%LET TYPE30DD = Y;** in USOURCE(GENGUIDE) to tell CPExpert that SMF Type 30(DD) information is available.
- ☐ If you wish CPExpert to analyze DASD performance from the perspective of critical workloads, take the following steps:
 - ☐ Define workload categories in USOURCE(DASGUIDE).
 - ☐ If you are running MVS Goal Mode), be sure that **%LET GOALMODE=Y;** has been specified in USOURCE(GENGUIDE).
 - ☐ Specify **%LET LOVED1 = xxxx** in USOURCE(DASGUIDE), where 'xxxx' is the name of a "loved one" workload.
 - ☐ Execute the DASCPE Module of the DASD Component.
- ☐ If you wish CPExpert to analyze DASD performance from the perspective of critical data sets, take the following steps **BEFORE** your performance data base has been updated:
 - ☐ Specify critical data set names and their associated response objectives in USOURCE(DASGUIDE).
 - ☐ Execute the DAS1415 Module to extract information from SMF Type 14/15 records.
- ☐ If you wish CPExpert to analyze DASD performance from the perspective of critical data sets, take the following steps **AFTER** your performance data base has been updated:
 - ☐ Make any desired modifications to data set names and their response objectives in USOURCE(DASGUIDE).
 - ☐ Specify **%LET DASDSN = Y** in USOURCE(DASGUIDE) to tell CPExpert to perform data set name analysis.
 - ☐ Execute the DASCPE Module of the DASD Component.

Section 5: DASD Analysis Factors

This section discusses DASD performance analysis considerations. Chapter 1 presents an overview of performance factors from a DASD I/O operation viewpoint. Chapter 2 highlights some of the factors that must be considered when analyzing DASD performance based upon data collected and recorded by SMF or RMF.

Chapter 1: Overview of DASD Performance Considerations

From a high-level view, there are four key measures of DASD performance: IOS Queue (IOSQ) time, pending (PEND) time, disconnect (DISC) time, and connect (CONN) time. These measures are reported by RMF in SMF Type 74 records. Exhibit 5-1 illustrates these four measures and another potential element of DASD I/O time, titled "Other".



MAJOR COMPONENTS OF DASD I/O OPERATIONS

EXHIBIT 5-1

Chapter 1.1: IOSQ time

IOSQ time is the time from the issuance of a STARTIO macro until the StartSubChannel (SSCH) instruction is issued. After the STARTIO macro is issued, the software determines whether the device is busy with *this system*; that is, whether there is an available Unit Control Block (UCB) for the device. If the device is not busy with *this system* (a UCB is available), the SSCH instruction is issued. However, if the device is busy with *this system*, the I/O request is queued. Thus, IOSQ time always means that the device is unable to handle additional requests from *this system*. (The emphasis on "this system" is explained in the below discussion of PEND time.)

This discussion of IOSQ time does not always apply to Parallel Access Volumes (PAVs)¹. With PAV devices, MVS creates multiple UCBs for each device, depending on how many “alias devices” have been defined. The multiple UCBs allow multiple active concurrent I/Os on a given device when the I/O requests originate from the same system². Using PAVs can dramatically improve I/O performance by nearly eliminating IOSQ.

Some small IOSQ time is often unavoidable. However, large IOSQ time imply a situation that should be examined. Large IOSQ times result from (1) too many I/O operations directed to the device or (2) lengthy device response times (perhaps caused by low percent cache hits or high PEND time). Large IOSQ times usually involve the following situations:

- Multiple data sets may be active on the volume. This situation is the most common and easiest to solve. The data sets can be redistributed among different volumes, to eliminate the queuing for the single volume. Alternatively, using the Parallel Access Volume feature available with IBM's Enterprise Storage Server (ESS) could allow multiple concurrent access to the device.
- Multiple users may be using the same data set on the volume. Depending upon the data set characteristics, duplicate copies of the data set placed on different volumes may solve the IOSQ problems. Alternatively, using the Parallel Access Volume feature available with IBM's Enterprise Storage Server (ESS) could allow multiple concurrent access to the device.
- Multiple application systems may be using the volume experiencing high IOSQ times. In this case, perhaps application redesign or scheduling can solve the problem. Alternatively, using the Parallel Access Volume feature available with IBM's Enterprise Storage Server (ESS) could allow multiple concurrent access to the device.
- A particular application (or system function) may be executing I/O to the device faster than the device can respond. Using application features as Data In Memory, increased buffering, using Local Shared Resources (LSR) or increasing buffer sizes, specifying optimal buffering parameters, and other similar enhancements could allow the applications to considerably reduce the use of I/O activity.
- The overall device response time (PEND, DISC, and CONN) times may be large, such that the device is unable to provide quick response to the I/O requests. This situation will be revealed by large values in the PEND, DISC, or CONN measures.

¹PAV devices are available with Enterprise Storage Server (ESS). With PAV devices, a “base device” address is defined, and a UCB is associated with this base address. “Alias device” addresses can be defined and UCBs are associated with the alias device addresses.

²Multiple Allegiance allows multiple active concurrent I/O operations on a given device when the I/O requests originate from different systems.

With Parallel Access Volumes (PAV) and dynamic alias management in Goal Mode, IOSQ time should be significantly reduced or eliminated. The implications of PAV and dynamic alias management will be discussed in rules related to these features.

Chapter 1.2: PEND time

PEND time is the time from the issuance of the StartSubChannel (SSCH) instruction until the device is selected by the control unit and physical positioning commands (such as seek and set sector) are transferred to the device. With modern fixed block architecture (FBA) devices, the PEND time ends when the physical positioning commands are presented to the *logical volume control block* within the control unit. The PEND time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, wait for device, or wait for “other” reasons)³.

The PEND time can be caused by the device being busy from *another system*. In this case, the system issuing the STARTIO macro (*this system*) would have no knowledge that the device was busy with another system. Rather, if a UCB were available for the device, the SSCH would be issued. However, the device could not necessarily be selected (unless multiple allegiance were available), since the device would be busy from another system.

Additionally, PEND time could accumulate even with PAV devices if the access were to an extent that was busy with another I/O operation from *this system*.

Large PEND times usually involve the following situations:

- **Shared devices.** If the device is shared with another system, PEND time may indicate contention with the other system. Large PEND times in shared-device environments usually involve situations very similar to those described under IOSQ time:
 - Multiple data sets may be active on the volume. This situation is the most common and easiest to solve. The data sets can be redistributed among different volumes, to eliminate the queuing at the channel level (reflected as PEND time) for the single volume.

Alternatively, if IBM's Enterprise Storage Server (ESS) is available, the Multiple Allegiance feature can be used to significantly reduce or eliminate PEND time caused by other systems. Multiple Allegiance allows multiple active concurrent I/O operations on a particular device when the I/O requests originate from different

³PEND time is significantly reduced with FICON channels. FICON channels can have multiple I/O operations concurrently active, which reduces the potential PEND time caused by channel busy. There is no port busy time with FICON switches, and control unit time is significantly reduced. This statement regarding PEND time is not necessarily correct if a large number (more than 5) I/O operations are concurrently executing on a FICON channel. Dr. H. Pat Artis and Mr. Robert Ross have presented the results of research indicating that performance degrades significantly when more than 5 I/O operations (Open Exchanges) are concurrently active on a FICON channel (see “Understanding FICON Channel Path Metrics” at www.perfassoc.com).

systems. With Multiple Allegiance, there is complete access with read I/O operations. For write I/O operations, there is concurrent access unless there is a conflicting extent⁴. If there is a conflicting extent, the controller holds the I/O operation in a PEND state for the device.

If some of the data sets are not required to be shared, then the Data Base Administrator has complete flexibility to move these data sets (subject, of course, to the performance implications of the target devices). These data sets could be moved to a non-shared device.

- Multiple applications or users may be using the same data set on the volume. Depending upon the data set characteristics, duplicate copies of the data set may be placed on different volumes. This would solve the PEND problems caused by contending systems. If this option is feasible, the data sets could be placed on non-shared devices, likely resulting in even more performance improvement. Alternatively, Record-Level Sharing (RLS) might provide a substantial reduction in the exclusive use of data sets.
- Multiple application systems may be using the volume experiencing high PEND times. In this case, perhaps application redesign or scheduling can solve the problem.
- **Non-shared devices.** Large PEND times for devices that are not shared may mean that there are insufficient paths available to the device. Too much I/O may be directed to many devices on the path, control unit, or. The data sets can be redistributed among different logical volumes on different paths, control units, or devices. This will reduce the hardware-level queuing. Alternatively, the entire volume may be moved to a different (less busy) path.

If redistributing the data sets or moving the logical volume is not feasible, then the device should have more paths. Depending upon the existing configuration, this may involve re-configuring existing channel paths, or acquiring additional hardware.

Fortunately, SMF Type 78 and Type 74 records contain information that can be used to identify at which level the hardware queuing occurs (that is, whether queuing is for the director port, control unit, or device; and CPEXpert calculates an estimate of the PEND time caused by channel activity).

- **Devices attached to cached controllers.** Large PEND times for devices attached to cached controllers may imply a high percent of read miss operations, or non-volatile storage (NVS) writes for IBM-3990-3 devices.

⁴ A conflicting extent is one in which the write operation attempts to update an extent.

To improve the probability of a read hit, the controller can *prestage* data into its cache. Prestaging means that data is read into the controller's cache ahead of its actually being required for use by an application. The amount of data that is prestaged depends on (1) whether the data is being accessed in a direct (random) mode or in a sequential mode and (2) the controller model and the enhancements made to the controller.

- C For *direct mode*, after the record is located, the 3390-3 and 3990-6 (initial version) stages in the balance of the track being read.

The 3990 Model 6 (with record cache) stages only the records requested into cache, eliminating the balance of the track staging that is normal with track caching as was implemented on initial versions of 3990-6 and on the 3990-3. This improvement reduces the PEND time caused by the controller busy during track staging.

- C As examples of prestaging for *sequential mode*, the 3990-3 reads up to two tracks into the cache⁵ before they are required, while the ESS 2105 sequential staging reads up to two cylinders ahead.

During prestaging operations for sequential reads, the control unit regularly checks to see whether other I/O requests are waiting to be processed. If any are waiting, the control unit interrupts the prestage operation, processes the queued requests, and continues with the prestage.

In DASD Fast Write Mode, the data is stored simultaneously in cache storage and in nonvolatile storage (NVS). At some subsequent time, the data in NVS can be *destaged* to DASD.

In Cache Fast Write Mode, data is placed into cache immediately, and there is no interaction with the device nor with NVS. However, if cache memory is required (or if Cache Fast Write Mode is turned off), the data in cache is destaged to DASD. Significant PEND time can result from destaging to DASD.

- **Dual Copy Initialize.** Large PEND times for IBM-3390 devices may be caused by dual copy initialize. IBM recommends the following⁶ for best system performance when using Dual Copy Initialize:

“C Use enhanced dual copy, that is, set DASD fast write on to all of your dual copy pairs.

⁵With the Sequential Staging Performance Enhancement, the 3990-3 can prestage up to a full cylinder (15 tracks) into the cache.

⁶Source: IBM 3993/9390 Introduction.

- C The write hit ratio should be 90% or higher for good performance. Write hit ratios are normally 99-100% if microcode supporting Record Cache II has been installed.
- C The read to write ratio should be 2:1 or greater for dual copy candidate volumes.
- C Wherever possible, use DLSE operations for your pairs.
- C As much as is possible, spread your dual copy pairs across multiple storage controls. By doing so, you lessen the impact of having a larger number of fast dual copy pairs on one subsystem, especially for a heavily loaded DASD subsystem. Remember that both devices in a duplex pair must be in the same logical DASD subsystem.”

Chapter 1.3: DISC time

DISC means that there is some delay that is often (but not always) associated with a mechanical movement during which the device disconnects from the control unit.

With legacy systems (e.g., 3380 drives attached to 3990-2 control units), the DISC time of most concern was associated with seek (arm movement) and rotational position sensing (time waiting for the disk platter to rotate to the location where desired data resides). Considerable performance improvement efforts were directed at reducing the seek activity and reducing the rotational position sensing (RPS)⁷ delays for the legacy systems. These two mechanical delays still exist for most modern *redundant array of independent disks* (RAID)⁸ systems, but their impact can not be directly reduced with normal methods.

With modern disks, data is cached into Actuator Level Buffers (ALBs), that contain data read from a track on the disk platter. Using ALBs can eliminate the RPS delays for records read on a particular track, since required data is read into the device buffer during a single rotation and stored until a path is available to transfer the data. However, if a record is to be read from a new track, some RPS delay could exist since the record would not be in the ALB, and must be read from the new track. Some initial RPS delay would apply in this case. This initial RPS delay is neither measured nor preventable.

Additionally, data is cached into increasingly large cache on the controller. For a read operation, desired data often is found in the cache. Write operations normally end as the data to be written is placed in non-volatile storage (NVS); and the storage processor writes

⁷ RPS delays are caused by a path not being available when the required data came under a device read head. Since a path was not available, the data could not be read and another rotation of the platter was experienced until the data again came under the device read head. Multiple rotations might be required, depending on the busy level of the path.

⁸ An array is an ordered collection of physical devices (disk drive modules) that are used to define logical volumes or devices.

the data to the device asynchronous with other activity (as a “back end” staging operation). The write activity can result in DISC time.

Consequently, DISC time for modern systems is a result of *cache read miss* operations, potentially back-end staging delay for write operations, peer-to-peer remote copy (PPRC) operations, and other miscellaneous reasons⁹. DISC time often can be very small with adequate cache. For example, there would be zero disconnect time for a cache read hit (the record was found in the cache).

Chapter 1.4: CONN time

CONN time includes the data transfer time, but also includes protocol exchange¹⁰ (or “hand shaking”) between the various components at several stages of the I/O operation.

For devices attached to paths that include parallel channels and ESCON channels, the data transfer time is simply the number of bytes transferred divided by the transfer speed. This is because a parallel channel or ESCON channel can have only one data transfer operation in execution at one time.

For devices attached to paths that include FICON channels, the algorithm is more complicated. This primarily is because a FICON channel can perform multiple data transfer (read and write) operations at one time. The data packets for multiple read or write operations are interleaved (or multiplexed) in the FICON link. CONN time for an individual I/O begins with the first frame of data transferred and ends last frame of data transfer, even though data for other I/O operations might be transferred concurrently on the link. Consequently, if multiple data packets (representing data for multiple read or write operations) are interleaved on the FICON link, the elapsed time for any particular I/O operation can be elongated¹¹ when compared with the elapsed time of the same I/O operation on an ESCON channel.

Chapter 1.5: OTHER time

There are at least two other potential I/O delays for DASD: (1) waiting for the I/O completion interrupt to be serviced by a processor and (2) waiting for the I/O interrupt to

⁹Artis has described a “sibling PEND” condition that results from collisions within the physical disk subsystem of RAID devices. See “Sibling PEND: Like a Wheel within a Wheel,” www.cmg.org/cmgpap/int449.pdf.

¹⁰Note that the protocol exchange occurs at multiple points in the normal I/O operation, even though it is shown only once in this exhibit.

¹¹The relative speed of a FICON channel is much higher than that of an ESCON channel. Consequently, the elapsed time of any particular I/O operation should be less on a FICON channel than on an ESCON channel, even if there are multiple I/O operations interleaving data. This statement regarding elapsed time is not necessarily correct if a large number (more than 5) I/O operations are concurrently executing on a FICON channel. Dr. H. Pat Artis and Mr. Robert Ross have presented the results of research indicating that performance degrades significantly when more than 5 I/O operations (Open Exchanges) are concurrently active on a FICON channel (see “Understanding FICON Channel Path Metrics” at www.perfassoc.com).

be serviced by a domain under PR/SM. Neither potential I/O delay is expected to be of the magnitude of the four "standard" I/O delays. However, they can be significant in special circumstances.

- Multi-processor configurations running under MVS can use any processor to service an I/O interrupt. However, when a processor services an I/O interrupt, the processor's high-speed cache storage is no longer valid when control is returned to the interrupted task. Consequently, many of the processor's high-performance design features may be nullified.

A hardware feature allows processors to be disabled for I/O interrupts. With this method, only a small number (perhaps only one) processor is enabled for interrupt processing. Only this processor will have its high-speed cache storage disturbed by the task-switching required for interrupt processing, and only this processor will periodically have its high-performance design features nullified. The disadvantage to this approach is that an interrupt may occur while the processor is busy servicing a previous interrupt.

If an interrupt is pending and no processor is enabled to service the interrupt, the interrupt must wait until a processor is available. This time should be insignificant, unless the system is processing a significantly large number of I/O operations. If the system is processing a large number of I/O operations, the interrupt pending delay could pose performance problems.

After the processor completes processing for an I/O interrupt, it issues a Test Pending Interrupt (TPI) instruction to determine whether there are any interrupts pending. If an I/O interrupt is pending, the processor proceeds to service that interrupt.

The IEAOPTxx member of SYS1.PARMLIB contains the **CPENABLE** keyword. This keyword specifies the percent of I/O interrupts detected by the TPI instruction, compared with all I/O interrupts. When the percent exceeds the high threshold of the CPENABLE keyword, MVS enables another processor to handle pending I/O interrupts. If the percent falls below the low threshold of the CPENABLE keyword, MVS will disable a processor (to the point that only one processor is enabled). The low and high threshold values for CPENABLE are 10 and 30 percent, respectively. These values normally mean that less than 30% of the I/O interrupts will be delayed for I/O interrupt service.

- MVS environments running under as a guest under VM or in a logical partition (LPAR) under PR/SM are subject to I/O interrupt delays. These delays can occur if another guest (for VM) or another domain is in its dispatch interval when the I/O interrupt completion is posted. The I/O interrupt remains pending until the guest or domain is dispatched. These delays have been estimated to be far more significant than might otherwise be expected.

Neither of the potential I/O delays described above is measured by RMF (although RMF does provide information on the number of I/O interrupts serviced by each processor and the number of TPI instructions resulting in I/O interrupt servicing).

The potential I/O delays are included in this discussion of general DASD performance considerations because (1) they may become important under certain situations and (2) techniques may be developed to assess their impact.

Chapter 2: RMF Data Analysis Considerations

This chapter highlights some of the factors that must be considered when analyzing DASD information collected and recorded by SMF in Type 30 records or by RMF in Type 70(series) records.

These factors do **not** preclude a comprehensive analysis of performance data and usually do not prevent insight into the causes of unacceptable performance. However, the factors must be recognized and accounted for both by CPExpert in analyzing data and by the user in reviewing CPExpert's results. The factors stem from (1) the way in which SMF and RMF create and record Type 30 and Type 70(series) information, and (2) inherent limitations caused by data averages.

Chapter 2.1: SMF information

SMF Type 30 records contain a record sub-type code to identify when the records are written:

CODE	SUB-TYPE DESCRIPTION
1	Job start
2	Interval records
3	Step termination
4	Step total
5	Job termination
6	System address space

SMF will optionally record the different sub-types, depending upon parameters contained in the SMFPRMxx member of SYS1.PARMLIB. Most installations collect Sub-type 4 (Step total) records, and many installations collect Sub-type 2 (Interval) records. If Interval records are recorded by SMF, Sub-type 3 (Step termination) records are automatically created.

It is highly desirable to collect Interval/Step termination information. It is virtually impossible to analyze system performance based upon Step total information if there exists long-running jobs. This is because it is impossible to correlate the information reflected in the Step total records with the information contained in SMF Type 70(series) data.

The Sub-type 4 records are written only after a long-running job step terminates, while the SMF Type 70(series) records are written at user-defined intervals (the interval typically is every 30 minutes or so). Long-running job steps may span many RMF recording intervals [RMF is responsible for creating the SMF Type 70(series) records]. Consequently, there

may be many RMF interval records written between the start and end of a long-running job step.

Sub-type 2 (Interval) records are written at user-defined intervals (typically the interval selected is the same interval as the RMF interval records). SMF writes a Sub-type 2 record when the specified interval has lapsed after the start of the job step and continues to write Sub-type 2 records at each subsequent interval. When the job step terminates, SMF writes a Sub-type 3 record containing the information since the last Sub-type 2 record was written. One consequence of the interval records is that system usage can be identified by workload, and can be correlated with the overall system statistics recorded by RMF in the SMF Type 70(series) records.

There are two variations in how SMF and RMF write interval data: (1) non-synchronized and (2) synchronized. Synchronization of SMF and RMF records is an option that must be explicitly specified in the SMFPRMxx member of SYS1.PARMLIB.

- C **Non-synchronized writing of SMF and RMF interval data.** With non-synchronized writing of SMF and RMF interval data, the Sub-type 2 records are written based upon the interval lapse from the start of the job step. They are not written at the same time as is the SMF Type 70(Series) records. This lack of coordination between recording the two record types poses a correlation problem: a particular Sub-type 2 (or Sub-type 3) record may span between two RMF recording intervals. From a data analysis view, there is no way to precisely determine whether the data reflected in the Sub-type 2 record (or Sub-type 3 record) should belong to the first RMF Type 70(series) interval or should belong to the second RMF Type 70(series) interval.

For example, suppose that the RMF recording interval were specified as 30 minutes, and RMF was directed to synchronize on the hour and half-hour. The RMF data would be collected and recorded at 10:00, 10:30, 11:00, 11:30, etc. Further suppose that a particular job step started at 15 minutes past the hour. Assuming that the Type 30 interval recording were specified as 30 minutes, SMF would create a Type 30 (Sub-type 2) interval record at 45 minutes past the hour, 15 past the next hour, and so forth. Thus, the RMF data would be recorded on the hour and half-hour, while the Sub-type 2 data would be recorded "offset" by 15 minutes.

CPEXpert addresses this problem by pro-rating the SMF Type 30 information based upon elapsed time. In the above example, 50% of the actual workload data contained in the SMF Type 30 (Sub-type 2 or Sub-type 3) records would be attributed to one RMF measurement interval and 50% would be attributed to the next RMF measurement interval. This pro-rating approach essentially assumes that the resources required by a job step do not vary much from one instant to the next.

This approach works quite nicely so long as the job step uses resources in a uniform fashion. Many job steps exhibit this characteristic, and the resources required by the job step do not vary much as the job executes. Resources distributed using the pro-

rating approach result in fairly consistent usage characteristics when comparing the summarized Type 30 data with SMF Type 70(series) data.

However, some job steps exhibit significant cycles, or require resources at the beginning or end of the job step. For these job steps, the pro-rating approach does not properly distribute the resource usage into the correct RMF measurement interval. Summarized Type 30 data would not compare well with Type 70(series) data if many job steps exhibit this cyclic or burst nature of resource usage. Unfortunately, there is no way to better distribute the data. Consequently, analysis based upon Type 30 data must be viewed with some caution. The analysis **generally** will be sufficiently precise for performance analysis purposes. However, anomalies will appear and results must always be subjected to a "reality" test.

This point is significant for the DASD Component, because the DASD Component attributes DASD usage to workloads based upon correlating SMF Type 30 data with SMF Type 74 data. The DASD I/O activity at the job step level is obtained from the SMF Type 30 interval records (using a modification to MXG or to MICS). This DASD I/O activity is pro-rated to the RMF measurement intervals as described above. The RMF DASD device I/O characteristics (IOSQ, PEND, DISC, and CONN times) are attributed to workloads based upon the pro-rating.

It is possible that the pro-rating method will result in improper attribution of I/O device characteristics to workloads. For example, suppose that a job step completed a few minutes past the hour (and that RMF data records were synchronized on the hour and half-hour). When the job step completed, it could execute many DASD I/O operations. These I/O operations would mostly be attributed to the previous RMF interval and the DASD device characteristics of that interval would be associated with most of the I/O operations. Suppose that there were no I/O problems with the device in the first interval. It is possible, however, that as the job step completed, it experienced significant DASD I/O problems that would be reflected in the second RMF interval. Since only a few of its I/O operations would be attributed to the second RMF interval, CPExpert would associate the I/O problems to only a few of its I/O operations. Consequently, CPExpert might consider the job step (and the workload category associated with the job step) to receive good DASD service, when the workload actually received bad service because of its burst I/O operations.

This example is not intended to invalidate the techniques CPExpert uses. Rather, the example is presented to explain a unique situation in which the techniques could result in improper conclusions. Readers may note that IBM's Service Level Reporter and any other software analyzing SMF/RMF data face the same analysis problem. The problem is with the data; not with the technique.

- C **Synchronized writing of SMF and RMF interval data.** Synchronized writing of SMF and RMF was introduced with MVS/ESA SP4. When interval accounting is synchronized, SMF generates interval records for a work unit based on the end of the

SMF global recording interval, rather than the start time of a job. This feature allows Type 30 records (and other record types) to be synchronized with writing RMF Type70(series) records. SMF places indicators (or "flags") in SMF Type70(series) records to indicate whether SMF and RMF records are synchronized.

It is not necessary for CPExpert to pro-rate data if the recording intervals are synchronized.

Chapter 2.2: Data Averages

The data collected by RMF and recorded in SMF Type 70(series) records provide a valuable source of information about the use and interaction of system components and workloads. However, the data are summarized and recorded at specific intervals (e.g., every 30 minutes). For most data elements, analysis must be accomplished based upon the **summary** or **average** values.

For example, the DASD IOSQ time reported for each device is the total for the measurement interval. The average IOSQ time per I/O operation is computed by dividing the total IOSQ time by the number of I/O operations. This average may have no relation to the IOSQ time experienced by any particular I/O operation. This problem is particularly pervasive as the RMF recording interval becomes more lengthy (e.g., if the recording interval were 60 minutes). The DASD IOSQ time may be quite long during the first half of the RMF measurement interval when contending workloads execute. The DASD IOSQ time may be short during the last half when a workload executes without contention from other applications. The average of the two extremes may lead to a conclusion that there was no problem with DASD IOSQ time for the entire interval!

Most DASD analysis performed by CPExpert assumes either a uniform or an exponential distribution of DASD I/O operations. For example, the pro-rating discussed in the previous chapter assumes a uniform distribution of I/O operations on a **job step** basis, over the life of the job step. However, the queuing models employed by CPExpert to analyze various aspects of DASD delays generally assume an exponential distribution of I/O operations at the **device** level, over an entire RMF measurement interval. Neither of these assumptions may be correct.

Wicks¹² illustrates a variety of distributions of the arrival rate of events, ranging from uniform distribution, to "cafeteria" distribution (events mostly arrive in clusters), to "London bus" distribution (events arrive only in clusters), to a random distribution (events exhibit a Poisson or exponential arrival). The arrival of many events in computer systems exhibit an exponential distribution, and M/M/1 or M/M/C queuing models can fairly represent many aspects of the systems.

¹²Wicks, R. J., *Balanced Systems and Capacity Planning*, IBM Corporation Washington Systems Center Technical Bulletin GG22-9299-02

However, Wicks gives an excellent example of exceptions: when editing a dataset using ISPF, the entire dataset may be read, some time is spent editing, and then the entire dataset may be written. The I/O requests in this instance would be similar to Wicks' "London Bus" distribution.

There is a tradeoff between (1) recording RMF data frequently, incurring the overhead and storage requirements of the additional RMF records, and requiring additional resources to analyze the data versus (2) recording RMF data less frequently, having less precise or representative data to analyze, and minimizing the resources required to perform the analysis. The importance of these tradeoffs must be evaluated in light of the objectives of the analysis and the requirements for precision of results.

In any event, any review of analysis and conclusions (whether by CPExpert or by a performance analyst) must be viewed with some caution because of the data summary, data averaging, and data distribution issues.

If the analysis consistently results in the same conclusions, you can be reasonably sure that the analysis is correct. However, it generally is unwise to make changes based upon analysis of a single day's RMF measurement information unless a "reality test" indicates that the analysis clearly is correct.

Section 6: Using the DASD Component

This section describes how to use the DASD Component to analyze system performance in the areas of DASD performance evaluated by the DASD Component.

As discussed in Section 1, the DASD Component is designed based upon the DASD performance analysis concept of "address the outrageous problems" rather than attempting to identify every potential problem. Most Data Base Administrators will be quite pleased to solve one major DASD problem at a time. Consequently, using the DASD Component involves (1) executing the software, (2) reviewing the output, (3) making system changes or altering the guidance to the DASD Component, and (4) iterating through the process.

This process is described in terms of providing guidance to the DASD Component, actions that should be taken on a daily basis (or more frequently if you wish), and actions that should be taken on a weekly or monthly basis.

Chapter 1: Prepare guidance for the DASD Component

Use TSO ISPF to change the CPEXPERT.USOURCE(DASGUIDE) PDS member to reflect the guidance required by the DASD Component. Section 3 describes the variables that must be changed in the DASGUIDE PDS member. This step must be taken only when the guidance changes (normally when you wish to select different to analyze, when you wish to exclude additional volumes, or when you wish to redefine workloads).

Chapter 2: Actions on a daily basis

Use the DASD Component to analyze overall system performance, by following the below steps:

Step 1: Execute the DASCPE Module

The JCL to execute the DASCPE Module is described in Exhibit 4-1. We suggest that you append the JCL to execute the DASCPE Module to the normal daily update of your performance data base.

Step 2: Review the output from the DASCPE Module

If any rules were produced, refer to the specific rule in Appendix A for a description of the rule, a discussion of why the rule was produced, and a recommendation for actions that should be taken.

Depending upon the output, you may wish to make changes or wait to see if the problems are identified in an analysis of a subsequent day's data.

- The DASD Component may identify problems which clearly should be solved because their effect is so serious. In many cases, once the problem is identified, users immediately realize that the problem and suggested solutions make sense.
- The DASD Component may identify problems which you do not feel will commonly occur. For example, you may review the "worst performing device" and realize that it is the "worst" only because of an infrequent application scheduling situation. In this case, you may wish to exclude the volume and reprocess the data. (The DASD Component will provide summary information about the next "worst" devices, so you can appreciate whether it is likely that there are other serious problems.)
- The DASD Component may identify problems about which you have doubts. The discussion in Section 5 illustrates the potential problems with the analysis, and the discussion associated with each rule often amplifies the cautions appropriate to the particular analysis. In this case, you probably should take no action but wait to see whether the same problems and analysis recur with subsequent data. If the same problems are identified by the DASD Component after analyzing several day's SMF/RMF data, you can be more confident that SMF/RMF data recording problems play less of a consideration.

Chapter 3: Actions on a weekly or monthly basis

Review the recommendations produced during the previous week or month. If the same problems consistently occur and CPExpert makes consistent recommendations, then you should consider action. This deliberate review of the problems and associated recommendations ensures that the problems are continuing and that change is warranted.

You generally should make only one change at a time! This sound tuning advice is founded on the principles that:

- Tuning is an art. No one (and certainly not CPExpert) can guarantee that any particular change will have a beneficial effect in all environments.
- Changes may have unexpected effects. Most systems are complex, parameters may improve performance of one area at the expense of performance in another area, and management may wish resources focused on the second area.
- If you make multiple changes and performance deteriorates, you will be unable to identify easily the change causing the problem. You are then faced with the problem of backing out all of the changes and starting over, one at a time.
- Some changes are not "precise" in that, for example, keyword values might need to be adjusted a little at a time until a suitable value is reached. If multiple changes are made, you will be unable to detect the effect of the fine-tuning of the changes.

Above all, **remember that the recommendations from CPExpert are simply options** to be considered in the context of overall objectives. You must decide whether the recommendations are reasonable. Rarely should a recommendation be implemented without first evaluating how the recommendation will effect other workloads.

Please remember that CPExpert is not intended to replace a performance analyst. Rather, CPExpert was developed to help analyze the performance of MVS systems. CPExpert automates much of the routine of computer performance evaluation. Performance analysts can then focus on the areas which are not routine and which "require thinking".

With this philosophy, please let us know when you discover areas in which CPExpert could have helped you analyze a problem. We will improve our product and you will have more help!

This request is particularly applicable to the DASD Component. We could have implemented much additional analysis and many additional options. We decided to wait to see what our users **wanted** and respond to your requirements, rather than trying to "guess" about your DASD analysis needs. However, we are eager to provide additional analysis; let us know what you want!

Description of Rules

Appendix A

This appendix contains a description of each rule that results in a finding by the DASD Component of CPEXpert. The description summarizes the rule, lists predecessor rules, discusses the rationale for the finding, and recommends action.

The summary of the rule presents a short description of the finding.

The predecessor rules are listed so you can follow the line of reasoning leading to a particular rule being executed.

The discussion describes as much as necessary of the operation of the computer system (the hardware, MVS, the Workload Manager, etc.) as it relates to the particular rule. The purpose of the discussion is to explain the reasoning behind the rule, and what causes the rule to be produced.

The recommendations suggest possible actions that should be considered based on the findings. In many cases, multiple possible actions are listed. You must determine which actions should be taken (this determination is based upon the suitability of the actions to your own environment, the financial implications of the action, and the "political" acceptability of the action.)

The rules are organized in numerical order. However, not all numbers are represented. The LIST OF RULES in this appendix lists all rules that are included in the initial release of the DASD Component.

The DAS2nn rules are very similar to the DAS1nn rules. The major difference in the DAS2nn rules is that they relate to the "expanded" analysis of "loved one" workload. Consequently, the narrative is somewhat different. However, the DAS2nn rules will often refer you to the DAS1nn rules for more detailed discussion and suggestions.

You may wish to read all of the rules in this appendix, just to see the type of problems that are encountered in different installations. However, it is not necessary to read all of the rules. It is necessary only to read the rules that apply to your installation. The rules that apply to your installation are identified by the report produced from the DASCPE Module.

Some of the rules (such as seek analysis) apply only to legacy systems (e.g., 3380 drives attached to 3990-2 control units). CPEXpert automatically detects the type of device and controller, and either invokes or suppresses rules that apply to legacy systems. The rules that apply only to legacy systems are marked with an asterisk in the listing of rules.

List of Rules

<u>RULE</u>	<u>DESCRIPTION</u>
DAS000	Sysplex performance characteristics of significant volumes
DAS050	Performance characteristics of significant volumes
DAS100	Volume with worst overall performance
DAS102	Volume with next worst overall performance
DAS105	Volume performance was not consistently poor in any area
DAS110*	Seeking was the major cause of I/O response delay
DAS111*	Seeking was probably caused by independent applications
DAS112*	Seeking was probably caused by a single application
DAS113*	Worst seeking was probably caused by independent applications
DAS114*	Worst seeking was probably caused by a single application
DAS115*	Seeking was cause of I/O delay on page pack
DAS120*	Missed RPS was major cause of I/O response delay
DAS121*	Volumes contributing to missed RPS
DAS123*	Non-DASD devices contributed to RPS delay
DAS125*	Applications contributing to RPS delay
DAS130	PEND time was major cause of I/O delay
DAS131	PEND delay time was caused by channel busy
DAS132	PEND delay time was caused by director port busy

* = These rules apply only to legacy systems.

<u>RULE</u>	<u>DESCRIPTION</u>
DAS133	PEND delay time was caused by controller busy delays
DAS134	PEND delay time was caused by device busy delays
DAS135	PEND time was caused by other delays
DAS140	High connect time was major cause of I/O response delays
DAS150	High IOSQ was major cause of I/O response delays
DAS160	Disconnect was major cause of response delay
DAS170	There did not appear to be a problem with the device
DAS180	Applications accessing the volume with the worst performance
DAS185	Applications accessing the volume during the period with worst performance
DAS200	Volume with worst overall performance from the perspective of the critical workload
DAS202	Volume with next worst overall performance from the perspective of the critical workload
DAS205	Volume performance was not consistently poor in any area
DAS210*	Seeking was the major cause of response delay to the critical applications
DAS220*	Missed rotational position sensing was major cause of response delay to the critical applications
DAS221*	Volumes contributing to missed rotational position sensing
DAS223*	Non-DASD devices contributed to RPS delay
DAS225*	Applications contributing to RPS delay

* = These rules apply only to legacy systems.

<u>RULE</u>	<u>DESCRIPTION</u>
DAS230	Large PEND time was major cause of response delay to the critical applications
DAS231	PEND delay time was caused by channel busy
DAS232	PEND delay time was caused by director port busy
DAS233	PEND delay time was caused by controller busy delays
DAS234	PEND delay time was caused by device busy delays
DAS235	PEND delay time was caused by other delays
DAS240	High connect time was major cause of response delays to the critical applications
DAS250	High IOSQ was major cause of response delays to the critical applications
DAS260	High disconnect time was major cause of response delay to the critical applications
DAS270	There did not appear to be a problem with the device
DAS280	Non-critical job steps used this volume and were a major cause of response delays to the critical applications
DAS285	Non-critical job steps used this volume during the period of worst performance and were a major cause of response delays to the critical applications
DAS287	Other applications did not reference the volume

* = These rules apply only to legacy systems.

<u>RULE</u>	<u>DESCRIPTION</u>
DAS300	Perhaps shared DASD caused performance problems
DAS385	Applications potentially causing worst shared DASD conflicts
DAS390	Shared DASD conflicts did not cause performance problems
DAS400	Access characteristics of significant data sets

<u>RULE</u>	<u>DESCRIPTION</u>
DAS600	Excessive Control Area (CA) splits occurred
DAS604	Excessive secondary extents were allocated
DAS605	Excessive extents were used and secondary allocation unit was small
DAS606	Primary or Secondary allocation unit was small
DAS607	VSAM data set is close to maximum number of extents
DAS610	Relatively small CI size was used for sequential processing
DAS611	Relatively large CI size was used for direct processing
DAS612	Relatively large CI size was used for mixed processing
DAS620	The number of data buffers should be increased
DAS621	The number of index buffers should be equal to index levels
DAS622	The number of index buffers should be more than STRNO value
DAS625	NSR was used, but a large percent of the access was direct
DAS635	LSR was used, but a large percent of the access was sequential

Rule DAS000: Sysplex performance characteristics of significant volumes

Finding: CPExpert identifies the performance characteristics of the volumes in the sysplex that have the most potential for performance improvement.

Impact: This finding is used to assess the importance of the "worst" performing device and to determine whether other devices offer significant performance improvement potential.

Logic flow: This is a basic finding. There are no predecessor rules.

Discussion: CPExpert uses the following algorithm to identify the devices that have the most potential for improvement:

- CPExpert computes the average device response time for each **type** of device in the configuration, for each RMF measurement interval. The logic computes the average device response by type of device, since better performance would be expected from cached devices (for example) than from non-cached devices. This method essentially assesses the performance of each device against the performance of similar devices in the configuration.
- Devices that exceed the average device response time for their device type in any RMF measurement interval are selected as candidates for improvement. The rationale is that improvement efforts should not be directed at devices that provide better than average response. Thus, the candidate set of devices to analyze consists of those that provided worse than average response.
- The I/O rate of each "candidate device is weighted by its response time, **for the entire set of RMF intervals in which the device exceeded the average response**. The result is a measure of the relative performance improvement **potential** of each device that provided worse than average response, from an overall system view. For example, consider two devices in a device type having an average I/O response of 20 milliseconds:

Device A: I/O rate = 30 I/O operations per second
 Device response = 25 milliseconds
 RMF intervals with above average response = 4
 Seconds per RMF interval = 900
 Weighting factor = $30 * 25 * 4 * 900 = 27,000,000$

Device B: I/O rate = 5 I/O operations per second
Device response = 40 milliseconds
RMF intervals with above average response = 5
Seconds per RMF interval = 900
Weighting factor = $5 * 40 = 900,000$

In the above example, CPExpert would select Device A as having the most overall potential for improvement, even though its per-I/O device response was not as bad as the device response of Device B.

CPExpert ranks the devices based on the weighting factor computed above. CPExpert then analyzes the devices, starting at the device with the highest weighting factor.

With Rule DAS000, CPExpert lists basic characteristics of the volumes having the most potential for improvement, so that you can appreciate the relative performance improvement potential between volumes on the list. The data presented by Rule DAS000 reflects the average per-second delays **only** during measurement intervals when the device I/O performance was worse than the average for its device type. This information is presented on a sysplex view basis, regardless of whether a specific system has been selected for analysis.

If the performance data base contains data for more than one sysplex, and if %LET SYSPLEX=*ALL; has been specified in USOURCE(DASGUIDE), CPExpert will produce information for all volumes in the performance data base. If a specific sysplex is selected for analysis (using the %LET SYSPLEX=xxxx, where "xxxx" is the name of a sysplex, only volumes for the designated sysplex will be listed.

The "weighted delays" value is a relative measure of the performance improvement potential of the volume. The absolute values in the column are not particularly meaningful. Rather, the values should be compared to each other to assess the relative performance impact of each volume.

It is possible that a volume may have a significant improvement potential in a particular measurement interval, but not be the volume with the most overall potential for improvement. This situation can arise because the analysis is directed toward the volumes with the **most overall** performance improvement potential. If you suspect that this is the case with a particular device, you can "select" that device for analysis, using the select process described in Section 3 of this document.

The following example illustrates the output from Rule DAS000:

RULE DAS000: SYSPLEX PERFORMANCE CHARACTERISTICS OF SIGNIFICANT VOLUMES

The following is a list of the most significant volumes showing their overall performance characteristics for the period being analyzed, from an overall sysplex view. The "average per second delays" represent the averages ONLY during measurement intervals when the device performance was worse than the average for this device type on at least one system in the sysplex. The "weighted delays" value is a measure of the overall relative performance impact of each volume.

SYSTEM	VOLSER	DEVICE NUMBER	I/O RATE	---	AVERAGE	PER	SECOND	DELAYS---	WEIGHTED DELAYS
				RESP	CONN	DISC	PEND	IOSQ	
SYZ0	SP0006	FE58	114.0	118.340	0.130	0.003	0.785	117.423	236679
SYH0	SP0006	FE58	93.2	94.382	0.106	0.003	0.775	93.498	188763
SYF0	SP0006	FE58	40.7	5.983	0.048	0.001	0.510	5.425	11967
SYE0	SP0006	FE58	28.5	3.444	0.035	0.001	0.404	3.005	6889
SY90	DB008D	FDFA	2.9	0.755	0.035	0.028	0.007	0.685	1509
SYA0	CAT00F	FD6B	1.5	0.413	0.018	0.001	0.006	0.387	826
SYF0	D83IA1	BDE1	7.0	0.319	0.195	0.000	0.024	0.099	638
SYA0	SP3057	FEC1	8.2	0.306	0.131	0.002	0.026	0.146	611
SYE0	CAT00F	FD6B	1.1	0.299	0.013	0.001	0.005	0.281	599
SYA0	DB0053	FEF9	1.4	0.298	0.014	0.003	0.004	0.277	597
SYA0	CAT011	FD7C	6.1	0.279	0.014	0.090	0.058	0.116	558
SYE0	CAT011	FD7C	5.3	0.256	0.013	0.065	0.061	0.117	512
SY80	CAT00F	FD6B	1.0	0.245	0.012	0.000	0.003	0.229	490

In this example, it is clear that SP0006 has significant performance improvement potential. The DASD Component would analyze SP0006 from the view of each system in which it appeared as the "worst" device, to determine what caused the delays. Additionally, if the CPExpert modification to MXG or MICS (described in Section 2) had been installed, the DASD Component would list the applications referencing SP0006. Further, if SMF Type 42 records were available (and the volume contained data sets managed by DFSMS), the DASD Component would produce Rule DAS400 to show access characteristics of the most significant data sets that resided on SP0006.

Notice that the data presented by Rule DAS000 are in "average per second" delays rather than "average per I/O" delays. This presentation gives the impact overall of each volume, which is appropriate for the weighted delays (or intensity) shown. If "average per I/O" delays were used, the effect of delays would not be as clear since devices with a few I/O operations could have significant delay per I/O operation. Displaying these significant delays would be misleading, since only a few I/O operations experienced the delays.

Suggestion: You should use the information displayed by Rule DAS000 to assess the relative importance of the "worst" performing device compared with the performance improvement potential of the other devices. In some cases, the impact of the "worst" performing device will be several times the impact of the next performing device. In most cases, the impact of the top five or six devices will account for most of the overall impact.

Rule DAS050: Performance characteristics of significant volumes

Finding: CPExpert identifies the performance characteristics of the volumes in a system that have the most potential for performance improvement.

Impact: This finding is used to assess the importance of the "worst" performing device and to determine whether other devices offer significant performance improvement potential.

Logic flow: This is a basic finding. There are no predecessor rules.

Discussion: CPExpert uses the following algorithm to identify the devices that have the most potential for improvement:

- CPExpert computes the average device response time for each **type** of device in the configuration, for each RMF measurement interval. The logic computes the average device response by type of device, since better performance would be expected from cached devices (for example) than from non-cached devices. This method essentially assesses the performance of each device against the performance of similar devices in the configuration.
- Devices that exceed the average device response time for their device type in any RMF measurement interval are selected as candidates for improvement. The rationale is that improvement efforts should not be directed at devices that provide better than average response. Thus, the candidate set of devices to analyze consists of those that provided worse than average response.
- The I/O rate of each "candidate device is weighted by its response time, **for the entire set of RMF intervals in which the device exceeded the average response**. The result is a measure of the relative performance improvement **potential** of each device that provided worse than average response, from an overall system view. For example, consider two devices in a device type having an average I/O response of 20 milliseconds:

Device A: I/O rate = 30 I/O operations per second
Device response = 25 milliseconds
RMF intervals with above average response = 4
Seconds per RMF interval = 900
Weighting factor = $30 * 25 * 4 * 900 = 27,000,000$

Device B: I/O rate = 5 I/O operations per second
Device response = 40 milliseconds
RMF intervals with above average response = 5
Seconds per RMF interval = 900
Weighting factor = $5 * 40 = 900,000$

In the above example, CPEXpert would select Device A as having the most overall potential for improvement, even though its per-I/O device response was not as bad as the device response of Device B.

CPEXpert ranks the devices based on the weighting factor computed above. CPEXpert then analyzes the devices, starting at the device with the highest weighting factor.

With Rule DAS050, CPEXpert lists basic characteristics of the volumes having the most potential for improvement, so that you can appreciate the relative performance improvement potential between volumes on the list. The data presented by Rule DAS050 reflects the average per-second delays **only** during measurement intervals when the device I/O performance was worse than the average for its device type. This information is presented on a system view basis.

The "weighted delays" value is a relative measure of the performance improvement potential of the volume. The absolute values in the column are not particularly meaningful. Rather, the values should be compared to each other to assess the relative performance impact of each volume.

It is possible that a volume may have a significant improvement potential in a particular measurement interval, but not be the volume with the most overall potential for improvement. This situation can arise because the analysis is directed toward the volumes with the **most overall** performance improvement potential. If you suspect that this is the case with a particular device, you can "select" that device for analysis, using the select process described in Section 3 of this document.

The following example illustrates the output from Rule DAS050:

RULE DAS050: PERFORMANCE CHARACTERISTICS OF SIGNIFICANT VOLUMES

The following is a list of the most significant volumes showing their overall performance characteristics for the period being analyzed. The "average per second delays" represent the averages ONLY during measurement intervals when the device I/O performance was worse than the average for this device type. The "weighted delays" value is a measure of the overall relative performance impact of each volume.

VOLSER	DEVICE NUMBER	I/O RATE	-----AVERAGE PER RESP	SECOND CONN	DELAYS----- DISC	PEND	IOSQ	WEIGHTED DELAYS
SY3085	72BF	77.6	0.282	0.112	0.002	0.030	0.138	27624
SVS10F	72FA	153.8	0.213	0.111	0.018	0.051	0.033	20879
DJ308D	3DD2	48.8	0.115	0.042	0.044	0.010	0.019	11303
SVCKC4	7054	3.4	0.113	0.027	0.000	0.086	0.000	11076
PS1345	3A21	118.5	0.112	0.074	0.004	0.023	0.011	8924
CFAC04	72CA	35.7	0.082	0.025	0.030	0.011	0.016	8029
EM300C	BB16	36.6	0.081	0.063	0.001	0.007	0.010	7916
DJ3012	3B81	39.5	0.078	0.028	0.033	0.008	0.009	7544
SY3062	7292	54.5	0.067	0.026	0.002	0.019	0.019	6529
SY3061	7291	12.0	0.055	0.035	0.005	0.005	0.010	5180
SVC102	7012	41.3	0.052	0.034	0.001	0.013	0.004	5131

In this example, it is clear that the top few devices have the most potential for improvement. The DASD Component would analyze SY3085 as the "worst" device, to determine what caused the delays. Additionally, if the CPExpert modification to MXG or MICS (described in Section 2) had been installed, the DASD Component would list the applications referencing SY3085. Further, if SMF Type 42 records were available (and the volume contained data sets managed by DFSMS), the DASD Component would produce Rule DAS400 to show access characteristics of the most significant data sets that resided on SY3085.

Notice that the data presented by Rule DAS050 are in "average per second" delays rather than "average per I/O" delays. This presentation gives the impact overall of each volume, which is appropriate for the weighted delays (or intensity) shown. If "average per I/O" delays were used, the effect of delays would not be as clear since devices with a few I/O operations could have significant delay per I/O operation. Displaying these significant delays would be misleading, since only a few I/O operations experienced the delays.

After analyzing SY3085, the DASD Component would analyze SVS10F as the next "worst" performing device, then analyze DJ308D, and so forth until the number of devices specified by the ANALYZE guidance variable in USOURCE(DASGUIDE) had been analyzed.

Suggestion: You should use the information displayed by Rule DAS050 to assess the relative impact of the "worst" performing device compared with the performance improvement potential of the other devices. In some cases,

the impact of the “worst” performing device will be several times the impact of the next performing device. In most cases, the impact of the top five or six devices will account for most of the overall impact.

Rule DAS055: Performance characteristics of significant volumes

Finding: CPExpert identifies the performance characteristics of the volumes in a system that have the most potential for performance improvement, **from the perspective of the “loved one” workload**.

Impact: This finding is used to assess the importance of the "worst" performing device from the perspective of the “loved one” workload, and to determine whether other devices offer significant performance improvement potential.

Logic flow: This is a basic finding. There are no predecessor rules.

Discussion: Rule DAS055 is similar to Rule DAS050, except that Rule DAS055 relates to devices accessed by “loved one” work. Please refer to Rule DAS050 for a discussion of the approach to selecting devices for analysis.

With Rule DAS055, CPExpert lists basic characteristics of the volumes having the most potential for improvement **from the perspective of the “loved one” workload**, so that you can appreciate the relative performance improvement potential between volumes on the list. The data presented by Rule DAS055 reflects the average per-second delays **only** during measurement intervals when the device I/O performance was worse than the average for its device type, **and** for measurement intervals when the device was referenced by the “loved one” workload. This information is presented on a system view basis.

The "weighted delays" value is a relative measure of the performance improvement potential of the volume. The absolute values in the column are not particularly meaningful. Rather, the values should be compared to each other to assess the relative performance impact of each volume.

It is possible that a volume may have a significant improvement potential in a particular measurement interval, but not be the volume with the most overall potential for improvement. This situation can arise because the analysis is directed toward the volumes with the **most overall** performance improvement potential. If you suspect that this is the case with a particular device, you can “select” that device for analysis, using the select process described in Section 3 of this document.

The following example illustrates the output from Rule DAS055 when a “loved one” workload was defined as BATCH:

RULE DAS055: PERFORMANCE CHARACTERISTICS OF SIGNIFICANT VOLUMES

The following is a list of the most significant volumes accessed by BATCH showing their overall performance characteristics for the period being analyzed. The "average per second delays" represent the averages ONLY during measurement intervals when the device I/O performance was worse than the average for this device type. The "weighted delays" value is a measure of the overall relative performance impact of each volume.

VOLSER	DEVICE NUMBER	I/O RATE	-----AVERAGE PER RESP	SECOND CONN	DELAYS----- DISC	PEND	IOSQ	WEIGHTED DELAYS
WORKP2	2324	32.0	0.320	0.305	0.005	0.008	0.001	8631
MVSPX1	232E	21.2	0.036	0.029	0.002	0.005	0.000	64
MVSPL1	2537	73.5	0.132	0.107	0.005	0.016	0.004	59
MVS902	2137	47.5	0.045	0.031	0.002	0.011	0.001	58
PRD015	2538	38.5	0.037	0.024	0.004	0.009	0.000	51
PRD002	252D	16.5	0.016	0.011	0.002	0.003	0.000	50
PRD004	2536	10.8	0.015	0.011	0.001	0.003	0.000	38
PRD017	2635	73.8	0.061	0.042	0.003	0.016	0.000	32
PRD007	2133	10.5	0.011	0.007	0.002	0.002	0.000	32
WORKP3	2220	7.4	0.072	0.065	0.004	0.002	0.002	30
PRD005	2038	24.8	0.027	0.018	0.002	0.006	0.000	29

In this example, WORKP2 has significant performance improvement potential from the perspective of the BATCH workload¹. The DASD Component would analyze WORKP2 as the "worst" device for the BATCH workload, to determine what caused the delays. Additionally, if the CPEXpert modification to MXG or MICS (described in Section 2) had been installed, the DASD Component would list the applications referencing WORKP2. Further, if SMF Type 42 records were available (and the volume contained data sets managed by DFSMS), the DASD Component would produce Rule DAS400 to show access characteristics of the most significant data sets that resided on WORKP2.

Suggestion: You should use the information displayed by Rule DAS055 to assess the relative impact of the "worst" performing device compared with the performance improvement potential of the other devices, from the perspective of the "loved one" workload.

In some cases (as shown in the above example), the impact of the "worst" performing device will be several times the impact of the next performing device. In most cases, the impact of the top five or six devices will account for most of the overall impact.

¹The BATCH workload was selected simply as an example of Rule DAS055.

Rule DAS100: VOLUME WITH WORST OVERALL PERFORMANCE

Finding: The identified volume had the worst overall performance during the entire measurement period. RULE DAS100 is similar to RULE DAS200; RULE DAS100 applies to **all** DASD devices, while RULE DAS200 applies only to DASD devices accessed by critical (or "loved one") workload.

Impact: The impact of this finding will depend upon the importance of the volume to overall system performance. If this is a critical volume, then this finding will have a HIGH impact. However, if the volume is accessed by low priority workloads, then this finding will have a LOW IMPACT or MEDIUM IMPACT.

Address spaces are retained in storage so long as they have uncompleted I/O operations. While the address spaces are in storage, they occupy page frames and may delay other address spaces. Additionally, the SRB time required to service I/O operations executes at a higher dispatching priority than a TCB, regardless of the dispatching priority of the TCB. Thus, there may be an overall system impact even though the volume may be accessed only by low priority workloads.

Logic flow: This is a basic rule finding; there are no predecessor rules.

Discussion: CPExpert determines the average device response time, by device type, for each measurement interval. A "device type" for this purpose is any unique device type (e.g., IBM-3380 or IBM-3390), with the device type modified to reflect whether the device is cached, is a Parallel Access Volume (PAV), or is a paging device.

The purpose of determining the average device response time, by device type, is the underlying principle that there is little point in analyzing a particular device if its response time is better than average. Rather, the most improvement potential resides with devices whose response time is worse than average.

CPExpert selects a device in each measurement interval for further analysis if the device response time exceeds the average for its device type.

CPExpert consolidates information in various SMF records to build a model of the I/O configuration. This model includes utilization and queuing information for all channel paths, controllers, and devices. In creating the model, CPExpert:

C Processes RMF Type 70 records to identify the systems that are in the sysplex.

C Processes RMF Type 73 records to identify the physical channels that are associated with each system, and the type of channel (e.g., ESCON, FICON-Bridge, FICON-Native, etc.). Additionally, the physical and LPAR channel busy time is acquired for each system.

C Processes RMF Type 78 records to obtain the logical control units associated with each system, and the channels associated with each logical control unit. Additionally, controller busy and director port busy times are acquired.

C Processes RMF Type 74 records to obtain devices associated with each logical control unit. Device performance characteristics are also acquired from the Type 74 records.

The result from the above is a record for each device, containing information about the devices; and the logical control units, channels, and systems that are associated with each device.

CPEXpert constructs a frequency distribution of all devices whose response is worse than average for the type of device, weighted by the number of I/O operations executed by the device. This yields a weighted measure of the potential performance improvement that might be achieved for each device. This frequency distribution is sorted descending, to yield an ordered list of the devices with the most improvement potential. This ordered list represents an "ordered intensity of access" distribution of the devices.

CPEXpert selects the top devices from the ordered list of devices with the most improvement potential. CPEXpert reports information about the top devices from the list, by sysplex and by system (see Rule DAS000 and Rule DAS050 for additional information about this information).

Detailed information regarding the "worst" devices is extracted and reported for each measurement interval.

For delays not directly measured (and contained in the RMF records), CPEXpert applies queuing formulae to the model to compute delays that occurred at significant parts of the model. The results from the model are associated with essential information describing the device response characteristics.

Exhibit DAS100-1 provides a sample output resulting from the analysis. The VOLSER and device number of the "worst" performing device are identified in the narrative. Information is provided about the overall average

I/O rate and the device utilization for the entire measurement period being analyzed.

RULE DAS100: VOLUME WITH WORST OVERALL PERFORMANCE

VOLSER SY3085 (device 72BF) had the worst overall performance during the entire measurement period (0:30, 31JUL2003 to 0:15, 01AUG2003). This pack had an overall average of 77.6 I/O operations per second, was busy processing I/O for an average of 14% of the time, and had I/O operations queued for an average of 14% of the time. Please note that percentages greater than 100% and Average Per Second Delays greater than 1 indicate that multiple I/O operations were concurrently delayed. This can happen, for example, if multiple I/O operations were queued or if multiple I/O operations were PENDING. The following summarizes significant performance characteristics of VOLSER SY3085:

MEASUREMENT INTERVAL	I/O RATE	--- AVERAGE PER SECOND DELAYS---	MAJOR PROBLEM
		RESP CONN DISC PEND IOSQ	
7:45- 8:00,31JUL2003	44.0	0.115 0.064 0.001 0.014 0.036	CONN TIME
8:00- 8:15,31JUL2003	41.1	0.137 0.058 0.001 0.014 0.064	QUEUEING
8:15- 8:30,31JUL2003	20.0	0.046 0.030 0.001 0.006 0.010	CONN TIME
8:30- 8:45,31JUL2003	35.1	0.094 0.050 0.001 0.012 0.031	CONN TIME
8:45- 9:00,31JUL2003	22.5	0.064 0.034 0.001 0.009 0.021	CONN TIME
9:00- 9:15,31JUL2003	38.6	0.107 0.055 0.001 0.014 0.037	CONN TIME
9:15- 9:30,31JUL2003	29.9	0.083 0.044 0.001 0.010 0.028	CONN TIME
9:30- 9:30,31JUL2003	1.8	0.004 0.003 0.000 0.001 0.000	
9:30- 9:45,31JUL2003	21.6	0.061 0.032 0.001 0.007 0.021	CONN TIME

VOLUME WITH WORST OVERALL PERFORMANCE

EXHIBIT DAS100-1

As shown in Exhibit DAS100-1, CPEXpert provides a summary for each measurement interval, showing the average I/O rate for the interval, and the **average delay time per second** during the interval (the total is shown as **I/O RESP** in Exhibit DAS100-1). The average delay time per second effectively reflects the percent of each second (shown in milliseconds) in which the associated delay occurred.

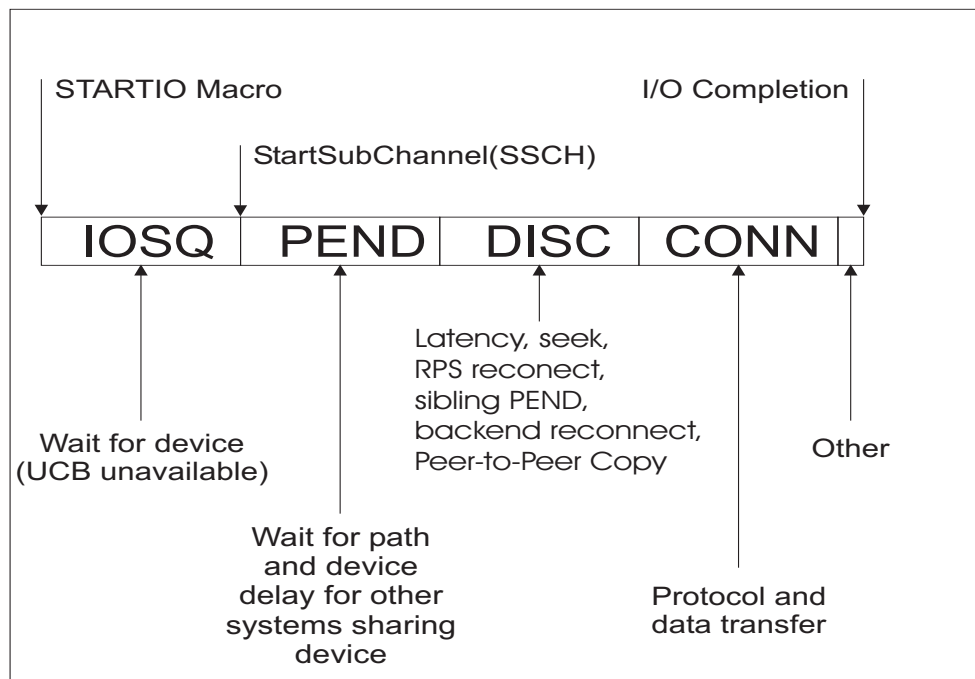
CPEXpert determines the major cause of device response delays (simply dividing each potential delay by the total device response time). The result from these calculations is evaluated to determine whether any area significantly predominates the I/O response time. If so, the respective area is listed as the major problem with the volume during the interval being analyzed.

The device may have no consistent problem. The problems may be concentrated in only a few measurement intervals and these few measurement intervals may dominate the performance characteristics of the volume. The worst of these intervals will be analyzed separately.

Additionally, it is possible that the device has no major problem for any interval. This condition most likely would occur if the guidance provided to CPEXpert is too restrictive. For example, if a large number of volumes are excluded from analysis (using the EXCLUDE option), the remaining volumes may have no particular problem. However, the logic is designed to select a "worst" volume, regardless of whether that volume actually has problems. CPEXpert tests for this condition and provides appropriate information.

From a high-level view, there are four key measures of DASD performance: IOS Queue (IOSQ) time, PENDING (PEND) time, disconnect (DISC) time, and connect (CONN) time. CPEXpert provides information on these four key measures, and identifies the major cause of response delay.

The following figure illustrates these four measures and another potential element of DASD I/O time, titled "Other":



C IOSQ time. IOSQ time is the time from the issuance of a STARTIO macro until the StartSubChannel (SSCH) instruction is issued. After the STARTIO macro is issued, the software determines whether the device is busy with *this system*; that is, whether there is an available Unit Control Block (UCB) for the device. If the device is not busy with *this system* (a UCB is available), the SSCH instruction is issued. However, if the device is busy with *this system*, the I/O request is queued. Thus, IOSQ time always means that the device is unable to handle additional requests from

this system. (The emphasis on "this system" is explained in the below discussion of PEND time.)

This discussion of IOSQ time does not always apply to Parallel Access Volumes (PAVs)¹. With PAV devices, MVS creates multiple UCBs for each device, depending on how many "alias devices" have been defined. The multiple UCBs allow multiple active concurrent I/Os on a given device when the I/O requests originate from the same system². Using PAVs can dramatically improve I/O performance by nearly eliminating IOSQ.

Please see Rule DAS150 for a more complete discussion of IOSQ time.

C PEND time. PEND time is the time from the issuance of the StartSubChannel (SSCH) instruction until the device is selected by the control unit and physical positioning commands (such as seek and set sector) are transferred to the device. With modern fixed block architecture (FBA) devices, the PEND time ends when the physical positioning commands are presented to the *logical volume control block* within the control unit. The PEND time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, wait for device, or wait for "other" reasons)³.

The PEND time can be caused by the device being busy from *another system*. In this case, the system issuing the STARTIO macro (*this system*) would have no knowledge that the device was busy with another system. Rather, if a UCB were available for the device, the SSCH would be issued. However, the device could not necessarily be selected (unless multiple allegiance were available), since the device would be busy from another system.

Additionally, PEND time could accumulate even with PAV devices if the access were to an extent that was busy with another I/O operation from *this system*.

Please see Rule DAS130 for a more complete discussion of PEND time.

¹PAV devices are available with Enterprise System Storage (ESS). With PAV devices, a "base device" address is defined, and a UCB is associated with this base address. "Alias device" addresses can be defined and UCBs are associated with the alias device addresses.

²Multiple Allegiance allows multiple active concurrent I/O operations on a given device when the I/O requests originate from different systems.

³PEND time is significantly reduced with FICON channels. FICON channels can have multiple I/O operations concurrently active, which reduces the potential PEND time caused by channel busy. There is no port busy time with FICON switches, and control unit time is significantly reduced. This statement regarding PEND time is not necessarily correct if a large number (more than 5) I/O operations are concurrently executing on a FICON channel. Dr. H. Pat Artis and Mr. Robert Ross have presented the results of research indicating that performance degrades significantly when more than 5 I/O operations (Open Exchanges) are concurrently active on a FICON channel (see "Understanding FICON Channel Path Metrics" at www.perfassoc.com).

C **DISC time.** DISC means that there is some delay that is often (but not always) associated with a mechanical movement during which the device disconnects from the control unit.

With legacy systems (e.g., 3380 drives attached to 3990-2 control units), the DISC time of most concern was associated with seek (arm movement) and rotational position sensing (time waiting for the disk platter to rotate to the location where desired data resides). Considerable performance improvement efforts were directed at reducing the seek activity and reducing the rotational position sensing (RPS)⁴ delays for the legacy systems. These two mechanical delays still exist for most modern *redundant array of independent disks* (RAID)⁵ systems, but their impact can not be directly reduced with normal methods.

With modern disks, data is cached into Actuator Level Buffers (ALBs), that contain data read from a track on the disk platter. Using ALBs can eliminate the RPS delays for records read on a particular track, since required data is read into the device buffer during a single rotation and stored until a path is available to transfer the data. However, if a record is to be read from a new track, some RPS delay could exist since the record would not be in the ALB, and must be read from the new track. Some initial RPS delay would apply in this case. This initial RPS delay is neither measured nor preventable.

Additionally, data is cached into increasingly large cache on the controller. For a read operation, desired data often is found in the cache. Write operations normally end as the data to be written is placed in non-volatile storage (NVS); and the storage processor writes the data to the device asynchronous with other activity (as a “back end” staging operation).

Consequently, DISC time for modern systems is a result of *cache read miss* operations, potentially back-end staging delay for write operations, peer-to-peer remote copy (PPRC) operations, and other miscellaneous reasons⁶. DISC time often can be very small with adequate cache. For example, there would be zero disconnect time for a cache read hit (the record was found in the cache).

Please see Rule DAS160 for a more complete discussion of DISC time.

⁴RPS delays are caused by a path not being available when the required data came under a device read head. Since a path was not available, the data could not be read and another rotation of the platter was experienced until the data again came under the device read head. Multiple rotations might be required, depending on the busy level of the path.

⁵An array is an ordered collection of physical devices (disk drive modules) that are used to define logical volumes or devices.

⁶Artis has described a “sibling PEND” condition that results from collisions within the physical disk subsystem of RAID devices. See “Sibling PEND: Like a Wheel within a Wheel,” www.cmg.org/cmgpap/int449.pdf.

-
- C **CONN time.** CONN time includes the data transfer time, but also includes protocol exchange⁷ (or "hand shaking") between the various components at several stages of the I/O operation.

For devices attached to paths that include parallel channels and ESCON channels, the data transfer time is simply the number of bytes transferred divided by the transfer speed. This is because a parallel channel or ESCON channel can have only one data transfer operation in execution at one time.

For devices attached to paths that include FICON channels, the algorithm is more complicated. This primarily is because a FICON channel can perform multiple data transfer (read and write) operations at one time. The data packets for multiple read or write operations are interleaved (or multiplexed) in the FICON link. CONN time for an individual I/O begins with the first frame of data transferred and ends last frame of data transfer, even though data for other I/O operations might be transferred concurrently on the link. Consequently, if multiple data packets (representing data for multiple read or write operations) are interleaved on the FICON link, the elapsed time for any particular I/O operation can be elongated⁸ when compared with the elapsed time of the same I/O operation on an ESCON channel.

Please see Rule DAS140 for a more complete discussion of CONN time.

- C **OTHER time.** There are at least two other potential I/O delays for DASD: (1) waiting for the I/O completion interrupt to be serviced by a processor and (2) waiting for the I/O interrupt to be serviced by a domain under PR/SM. Neither potential I/O delay is expected to be of the magnitude of the four "standard" I/O delays. However, they can be significant in special circumstances.
- C Multi-processor configurations can use any processor to service an I/O interrupt. However, when a processor services an I/O interrupt, the processor's high-speed cache storage is no longer valid when control is returned to the interrupted task. Consequently, many of the processor's high-performance design features may be nullified.

⁷Note that the protocol exchange occurs at multiple points in the normal I/O operation, even though it is shown only once in this exhibit.

⁸The relative speed of a FICON channel is much higher than that of an ESCON channel. Consequently, the elapsed time of any particular I/O operation should be less on a FICON channel than on an ESCON channel, even if there are multiple I/O operations interleaving data. This statement regarding elapsed time is not necessarily correct if a large number (more than 5) I/O operations are concurrently executing on a FICON channel. Dr. H. Pat Artis and Mr. Robert Ross have presented the results of research indicating that performance degrades significantly when more than 5 I/O operations (Open Exchanges) are concurrently active on a FICON channel (see "Understanding FICON Channel Path Metrics" at www.perfassoc.com).

A hardware feature allows processors to be disabled for I/O interrupts. With this method, only a small number (perhaps only one) processor is enabled for interrupt processing. Only this processor will have its high-speed cache storage disturbed by the task-switching required for interrupt processing, and only this processor will periodically have its high-performance design features nullified. The disadvantage to this approach is that an interrupt may occur while the processor is busy servicing a previous interrupt.

If an interrupt is pending and no processor is enabled to service the interrupt, the interrupt must wait until a processor is available. This time should be insignificant, unless the system is processing a significantly large number of I/O operations. If the system is processing a large number of I/O operations (or if the I/O is particularly time-sensitive), the interrupt pending delay could pose performance problems.

After the processor completes processing for an I/O interrupt, it issues a Test Pending Interrupt (TPI) instruction to determine whether there are any interrupts pending. If an I/O interrupt is pending, the processor proceeds to service that interrupt.

The CPENABLE keyword in the IEAOPTxx member of SYS1.PARMLIB is used to specify the percent of I/O interrupts detected by the TPI instruction, compared with all I/O interrupts. When the percent exceeds the high threshold of the CPENABLE keyword, MVS enables another processor to handle pending I/O interrupts. If the percent falls below the low threshold of the CPENABLE keyword, MVS will disable a processor (to the point that only one processor is enabled). IBM's recommended setting for the CPENABLE keyword differs, depending on the level of processor.

- C MVS environments running under as a guest under VM or in a logical partition (LPAR) under PR/SM are subject to I/O interrupt delays. These delays can occur if another guest (for VM) or another domain is in its dispatch interval when the I/O interrupt completion is posted. The I/O interrupt remains pending until the guest or domain is dispatched. These delays have been estimated to be far more significant than might otherwise be expected.

Suggestion: There are no suggestions directly associated with this rule. Subsequent rules will analyze the device problems and attempt to determine the cause of poor performance.

Rule DAS102: VOLUME WITH NEXT WORST OVERALL PERFORMANCE

Finding: The identified volume had the next worst overall performance during the entire measurement period.

Impact: The impact of this finding will depend upon the importance of the volume to overall system performance. If this is a critical volume, then this finding will have a HIGH impact. However, if the volume is accessed by low priority workloads, then this finding will have a LOW IMPACT or MEDIUM IMPACT.

Address spaces are retained in storage so long as they have uncompleted I/O operations. While the address spaces are in storage, they occupy page frames and may delay other address spaces. Additionally, the SRB time required to service I/O operations executes at a higher dispatching priority than a TCB, regardless of the dispatching priority of the TCB. Thus, there may be an overall system impact even though the volume may be accessed only by low priority workloads.

Logic flow: This is a basic rule finding; there are no predecessor rules.

Discussion: Rule DAS100 identified the volume that had the worst overall performance during the entire measurement period. Rule DAS102 is produced for each successive "worst performing" device selected from an ordered list (Rule DAS050 shows the ordered list of devices).

The number of devices analyzed by Rule DAS102 (and successive rules resulting from the analysis of each device) is controlled by the **ANALYZE** guidance variable (see Section 3: Specifying Guidance Variables).

Exhibit DAS102-1 provides a sample output resulting from the analysis. The VOLSER and device number of the "worst" performing device are identified in the narrative. Information is provided about the overall average I/O rate and the device utilization for the entire measurement period being analyzed.

Please refer to Rule DAS100 for a discussion of the information presented with Rule DAS102.

Suggestion: There are no suggestions directly associated with this rule. Subsequent rules will analyze the device problems and attempt to determine the cause of poor performance.

RULE DAS102: VOLUME WITH NEXT WORST OVERALL PERFORMANCE

VOLSER RSA002 (device 72FA) had the next worst overall performance during the entire measurement period (0:30, 31JUL2003 to 0:15, 01AUG2003). This pack had an overall average of 153.8 I/O operations per second, was busy processing I/O for an average of 18% of the time, and had I/O operations queued for an average of 3% of the time. Please note that percentages greater than 100% and Average Per Second Delays greater than 1 indicate that multiple I/O operations were concurrently delayed. This can happen, for example, if multiple I/O operations were queued or if multiple I/O operations were PENDING. The following summarizes significant performance characteristics of VOLSER SVS10F:

MEASUREMENT INTERVAL	I/O RATE	---	AVERAGE RESP	PER SECOND CONN	DELAYS---	DISC	PEND	IOSQ	MAJOR PROBLEM
0:30- 0:45,31JUL2003	88.2	0.107	0.063	0.012	0.028	0.004	0.000	CONN TIME	
0:45- 1:00,31JUL2003	26.4	0.033	0.020	0.004	0.009	0.000	0.000	CONN TIME	
1:03- 1:15,31JUL2003	23.2	0.027	0.017	0.003	0.007	0.000	0.000	CONN TIME	
1:15- 1:30,31JUL2003	134.0	0.315	0.095	0.024	0.052	0.144	0.000	QUEUEING	
1:30- 1:45,31JUL2003	90.9	0.130	0.065	0.012	0.037	0.016	0.000	CONN TIME	

VOLUME WITH WORST OVERALL PERFORMANCE

EXHIBIT DAS102-1

Rule DAS105: VOLUME PERFORMANCE WAS NOT CONSISTENTLY POOR

Finding: The performance of the volume having the "worst" overall performance was not consistently poor in any single area.

Impact: This finding has little impact, other than to indicate that there may be no problem with the device.

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with worst overall performance

Discussion: CPExpert checked PEND time, CONN time, DISC time, and IOSQ time. For legacy devices, CPExpert also has checked seek time and missed RPS reconnect time. None of these areas **consistently** accounted for a majority of the delay for the volume selected as the volume with the "worst" overall performance.

Suggestion: This condition most likely would occur when the period(s) of poor performance were extremely poor or if the guidance provided to CPExpert is too restrictive.

- Devices sometimes have short periods of extremely poor performance combined with a large number of I/O operations. This combination may result in the device being selected as the one with the most potential for performance improvement, even though the device did not have consistently poor performance. The most likely cause of such a situation is a particular application or combination of applications executing in the interval with poor performance.
- The guidance provided to CPExpert may be too restrictive. For example, if a large number of volumes are excluded from analysis (using the EXCLUDE option), the remaining volumes may have no particular problem. However, the logic is designed to select a "worst" volume, regardless of whether that volume actually has problems.

Rule DAS110: SEEKING WAS THE MAJOR CAUSE OF RESPONSE DELAY

Finding: Seeking was the major cause of the I/O response delay with the device.

Impact: This finding can have a MEDIUM IMPACT or HIGH IMPACT, depending upon the importance of the applications referencing the device and on the amount of seeking being done. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).*

Logic flow: The following rule causes this rule to be invoked:
 DAS100: Volume with the worst overall performance

Discussion: As described in Chapter 1 of Section 5 of this User Manual, DISC time with legacy systems is composed of seeking, latency, and missed RPS reconnect. The discussion associated with RULE DAS100 describes how CPExpert consolidates information in various SMF records to build a model of the I/O configuration. This model includes utilization and queuing information for all channel paths, controllers, and devices.

If legacy systems are in the configuration, CPExpert applies queuing formulae to the model to estimate the amount of delay attributed to missed RPS reconnect (these delays are a function of the probability of a device finding all paths busy when the device tries to reconnect to the channel path).

The estimated missed RPS reconnect time is subtracted from the DISC time reported by RMF for the device. Additionally, the average latency for the device type is subtracted from the DISC time. The resulting time is assumed to be the seek time. (Note the below discussion about why this assumption might not be correct.)

CPExpert performs the above analysis for each measurement interval reflected in the data. RULE DAS110 is produced if seeking was the major problem for a majority of the measurement intervals.

There are potential problems with this approach, although the approach is generally used throughout the computer industry as a way of estimating missed RPS delays and of estimating seeking for legacy devices.

- The queueing formulae assume exponential interarrival times, exponential service distributions, and an infinite population (the M/M/c formula - Erlang's C formula - is used for the calculations). These assumptions

may not be correct if, for example, the I/O activity is a function of a single application.

In his class "MVS I/O Configuration Management", Dr. Jeffery Buzen provides a Dump/Restore application as an excellent example of an application that does not follow standard queuing assumptions.

- The device may be cached, and it may be impossible to apportion the DISC time residual after subtracting missed RPS reconnect time. This time may represent a few missed cache read operations with long seek distances, or may represent a relatively large number missed cache read operations with little seeking but the standard latency for the device.

Thus, the seeking analysis can only show potential problems, rather be considered a definitive indication. However, it is usually a fairly accurate indication of the problem. If high average seeking is reported, you can be fairly certain that high seeking did occur. This is particularly true if the problem is reported throughout the measurement intervals. The "uncertainty" tends to be related to relatively low seeking or seeking reported for cached devices.

RULE DAS110 reports the overall average number of milliseconds out of each second in which the device was positioning the arm. Additionally, RULE DAS110 summarizes key information about the period of worst performance, if seeking was the major cause of delay during this period.

Suggestion: The seeks can be minimized by (1) rearranging files within the pack, (2) moving files from the pack to another actuator, (3) changing the application file accessing characteristics, or (4) possibly restricting the applications allowed to access the pack.

Rule DAS111: SEEKING WAS PROBABLY CAUSED BY INDEPENDENT APPLICATIONS

Finding: CPExpert determined that seeking was the major cause of delay in DASD response for the device. More than half of the I/O queuing to the device was explained using a queuing model. Consequently, CPExpert believes that the seeking probably was caused by independent applications.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device. If the finding is correct, and the independent applications are referencing different files, then the corrective actions could result in significant performance improvements. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).* |

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance
DAS110: Seeking was the major cause of response delay

Discussion: CPExpert uses a M/M/1 queuing model to calculate an estimated queue time for each measurement interval being analyzed. The underlying assumptions of the model are exponential interarrival times, exponential service distributions, and an infinite population. If device activity occurs in this way, the queuing model can predict the expected queuing delays.

If the queuing delay as measured by RMF is significantly different from the estimated delay from the model, it would be clear that the activity did not occur in a random fashion, and most likely the cause of the difference would be that the interarrival times are not randomly distributed.

However, the I/O delay to this pack was fairly well explained by the queuing model. More than half of the I/O queuing delay was explained by the model, for a majority of the measurement intervals. This indicates that the activity was mostly random, as would be expected from independent applications accessing files on the pack.

Suggestion: The most improvement for this pack would likely result from (1) separating the files to different packs, (2) rearranging the files within the pack, (3) tuning the file structure (for example, compressing a shared partitioned data set (PDS), (4) scheduling the applications to avoid contention, or (5) examining the applications doing most of the I/O.

You can identify the files with the most activity in one of several ways, depending upon the volume and the application involved. The method you select should depend upon the availability of information and tools in your environment.

- Discuss the file activity with applications personnel or systems personnel (depending upon the volume). If knowledgeable personnel are available, this is often the quickest and easiest way to determine how to solve the problems. Once a file access problem has been brought to their attention, applications or systems personnel often can easily decide upon a good solution.
- Use an exit available in MXG or in MICS to select Type 30(DD) information just for the volume. The Type 30(DD) information is rarely retained in a performance data base because it is so voluminous. However, you easily can code an exit in MXG or MICS to select the Type 30(DD) information for a specific volume. The amount of data selected would not be too large in most cases.

You can then write a SAS program to list the DD names used by applications, weighted by the number of accesses. For example, you could code the following to analyze data extracted during MXG update processing:

```
PROC FREQ DATA=pdblib.file;  
  WHERE DEVNR = addressX;  
  TABLES DDNAME/OUT=TEMP NOPRINT;  
  WEIGHT EXCPS;  
PROC SORT DATA=TEMP;  
  BY DESCENDING COUNT;  
PROC PRINT;  
RUN;
```

The device address is displayed by RULE DAS100, so the "address" in the above coding would be replaced with the actual address displayed by that rule. If you have a MXG performance data base, the address is retained in hexadecimal format, so you would suffix an X to the address.

If you have a MICS performance data base, the variable names must be changed appropriately (for example, DEVNR would be changed to DEVADDR). Additionally, the address is retained in MICS as a character representation of the value, so you do not need to suffix an X to the address. Rather, you must enclose the address in quotes.

The result from the above code would be a list of all DDNAMEs referencing the device being analyzed, weighted by the number of

EXCPs to the device, and ordered descendingly with the DD statements for the most active files listed first. You often cannot be certain that the most referenced files are causing problems. However, in most cases, you will find that only a few files (often only 2 or 3) account for over 90% of the accesses. These files generally will be the ones causing arm contention problems.

You can then write a SAS program to list the applications (jobs) referencing the device, weighted by the number of accesses. For example, you could code the following to analyze data extracted during MXG update processing:

```
PROC FREQ DATA=pdblib.file;  
  WHERE DEVNR = addressX;  
  TABLES JOB/OUT=TEMP NOPRINT;  
  WEIGHT EXCPS;  
PROC SORT DATA=TEMP;  
  BY DESCENDING COUNT;  
PROC PRINT;  
RUN;
```

The result from the above code would be a list of all jobs referencing the device being analyzed, weighted by the number of EXCPs to the device, and ordered descendingly with the jobs with the most active files on the device listed first. These jobs generally will be the ones causing arm contention problems. (If you have a MICS performance data base, you would change the code as described earlier.)

- Use a commercially-available DASD activity monitor to isolate the files accessed on the problem volume. If such a monitor is available, this would be a direct way to determine the problems.

Rule DAS112: SEEKING WAS PROBABLY CAUSED BY A SINGLE APPLICATION

Finding: CPEXpert determined that seeking was the major cause of delay in DASD response for the device. Less than half of the I/O queuing to the device was explained using a queuing model. Consequently, CPEXpert believes that the seeking probably was caused by a single application, rather than by independent applications.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device. If the finding is correct, then actions directed to the specific application could result in significant performance improvements. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).*

Logic flow: The following rules cause this rule to be invoked:
 DAS100: Volume with the worst overall performance
 DAS110: Seeking was the major cause of response delay

Discussion: CPEXpert uses a M/M/1 queuing model to calculate an estimated queue time for each measurement interval being analyzed. The underlying assumptions of the model are exponential interarrival times, exponential service distributions, and an infinite population. If device activity occurs in this way, the queuing model can predict the expected queuing delays.

If the queuing delay as measured by RMF is significantly different from the estimated delay from the model, it would be clear that the activity did not occur in a random fashion, and most likely the cause of the difference would be that the interarrival times are not randomly distributed.

This was the case with the volume being analyzed. Consequently, CPEXpert concludes that the activity was caused by a single application accessing the device. Note that there could be more than one application accessing the device, but if their access patterns were not random with respect to each other, then this conclusion would still be valid.

Suggestion: The most improvement for this pack would likely result from (1) separating the files to different packs, (2) rearranging the files within the pack, (3) tuning the file structure (for example, compressing a shared partitioned data set (PDS), or (4) examining the application doing most of the I/O.

You can identify the files with the most activity in one of several ways,

depending upon the volume and the application involved. The method you select should depend upon the availability of information and tools in your environment.

- Discuss the file activity with applications personnel or systems personnel (depending upon the volume). If knowledgeable personnel are available, this is often the quickest and easiest way to determine how to solve the problems. Once a file access problem has been brought to their attention, applications or systems personnel often can easily decide upon a good solution.
- Use an exit available in MXG or in MICS to select Type 30(DD) information just for the volume. The Type 30(DD) information is rarely retained in a performance data base because it is so voluminous. However, you easily can code an exit in MXG or MICS to select the Type 30(DD) information for a specific volume. The amount of data selected would not be too large in most cases.

You can then write a SAS program to list the DD names used by applications, weighted by the number of accesses. For example, you could code the following to analyze data extracted during MXG update processing:

```
PROC FREQ DATA=pdblib.file;  
  WHERE DEVNR = addressX;  
  TABLES DDNAME/OUT=TEMP NOPRINT;  
  WEIGHT EXCPS;  
PROC SORT DATA=TEMP;  
  BY DESCENDING COUNT;  
PROC PRINT;  
RUN;
```

The device address is displayed by RULE DAS100, so the "address" in the above coding would be replaced with the actual address. If you have a MXG performance data base, the address is retained in hexadecimal format, so you should suffix an X to the address.

If you have a MICS performance data base, the address is retained as a character representation of the value, so you do not need to suffix an X to the address. Rather, you must enclose the address in quotes.

The result from the above code would be a list of all DDNAMEs referencing the device being analyzed, weighted by the number of EXCPs to the device, and ordered descendingly with the DD statements for the most active files listed first. You often cannot be certain that the most referenced files are causing problems. However,

in most cases, you will find that only a few files (often only 2 or 3) account for over 90% of the accesses. These files generally will be the ones causing arm contention problems.

- Use a commercially-available DASD activity monitor to isolate the files accessed on the problem volume. If such a monitor is available, this would be a direct way to determine the problems.

Rule DAS113: WORST SEEKING WAS PROBABLY CAUSED BY INDEPENDENT APPLICATIONS

Finding: CPEXpert determined that seeking was the major cause of delay in DASD response for the device during the measurement interval with the worst I/O response. More than half of the I/O queuing to the device was explained using a queuing model. Consequently, CPEXpert believes that the seeking during this interval probably was caused by independent applications.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device. If the finding is correct, and the independent applications are referencing different files, then the corrective actions could result in significant performance improvements. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).* |

Logic flow: The following rules cause this rule to be invoked:
 DAS100: Volume with the worst overall performance
 DAS110: Seeking was the major cause of response delay

Discussion: CPEXpert uses a M/M/1 queuing model to calculate an estimated queue time for the measurement interval with the worst I/O response for the device being analyzed. The underlying assumptions of the model are exponential interarrival times, exponential service distributions, and an infinite population. If device activity occurs in this way, the queuing model can predict the expected queuing delays.

If the queuing delay as measured by RMF is significantly different from the estimated delay from the model, it would be clear that the activity did not occur in a random fashion, and most likely the cause of the difference would be that the interarrival times are not randomly distributed.

However, the I/O delay to this pack was fairly well explained by the queuing model. More than half of the I/O queuing delay was explained by the model, for a majority of the measurement intervals. This indicates that the activity was mostly random, as would be expected from independent applications accessing files on the pack.

Suggestion: The most improvement for this pack would likely result from (1) separating the files to different packs, (2) rearranging the files within the pack, (3) tuning the file structure (for example, compressing a shared partitioned data set (PDS), (4) scheduling the applications to avoid contention, or (5) examining the applications doing most of the I/O.

The actions to be taken when this rule is produced are the same associated with RULE DAS111. Please refer to that rule for further suggestions.

Rule DAS114: WORST SEEKING WAS PROBABLY CAUSED BY A SINGLE APPLICATION

Finding: CPExpert determined that seeking was the major cause of delay in DASD response for the device, during the measurement interval with the worst I/O response. Less than half of the I/O queuing to the device was explained using a queuing model. Consequently, CPExpert believes that the seeking probably was caused by a single application, rather than by independent applications.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device. If the finding is correct, then actions directed to the specific application could result in significant performance improvements. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).*

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance
DAS110: Seeking was the major cause of response delay

Discussion: CPExpert uses a M/M/1 queuing model to calculate an estimated queue time for the measurement interval with the worst I/O response for the device being analyzed. The underlying assumptions of the model are exponential interarrival times, exponential service distributions, and an infinite population. If device activity occurs in this way, the queuing model can predict the expected queuing delays.

If the queuing delay as measured by RMF is significantly different from the estimated delay from the model, it would be clear that the activity did not occur in a random fashion, and most likely the cause of the difference would be that the interarrival times are not randomly distributed.

This was the case with the volume being analyzed. Consequently, CPExpert concludes that the activity was caused by a single application accessing the device. Note that there could be more than one application accessing the device, but if their access patterns were not random with respect to each other, then this conclusion would still be valid.

Suggestion: The most improvement for this pack would likely result from (1) separating the files to different packs, (2) rearranging the files within the pack, (3) tuning the file structure (for example, compressing a shared partitioned data set (PDS), or (4) examining the application doing most of the I/O.

The actions to be taken when this rule is produced are the same associated with RULE DAS112. Please refer to that rule for further suggestions.

Rule DAS115: SEEKING WAS CAUSE OF I/O DELAY ON PAGE PACK

Finding: CPExpert determined that seeking was the major cause of delay in DASD response for the device with the worst performance, and the device was a paging volume.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device. High seeking on paging volumes potentially can have a significant impact on overall system performance if overall page delay is too high. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).*

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance
DAS110: Seeking was the major cause of response delay

Discussion: CPExpert has determined that seeking was the major cause of I/O response delays on the volume with the worst performance. CPExpert noticed that the device was a page pack.

Suggestion: If this finding resulted from an analysis of your entire DASD environment, then the finding is quite significant.

However, if you had requested an analysis of only the paging volumes (using the "select volumes for analysis" option in DASGUIDE), then the finding may not be significant. Please keep in mind that the logic of the DASD Component will always identify a "worst" device.

You should use the MVS Component or the TSO Component to assess whether paging delays are a major cause of performance degradation.

- If paging delays are not a major cause of performance degradation, then this finding probably should be ignored.
- If paging delays are a major cause of performance degradation, then you should consider adding another page pack to the configuration for this system.

Rule DAS120: MAJOR CAUSE OF I/O DELAY WAS MISSED RPS RECONNECT

Finding: CPExpert has determined that missed Rotational Position Sensing (RPS) reconnects was a major cause of delay in DASD response for the device.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).*

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance

Discussion: CPExpert computes the average channel path busy for paths to the device. SMF Type 78 information is used to acquire the channel path utilization, and the IOCP macro information is used to determine the channel paths to the device.

If the device is capable of DLS mode, two channel paths can be concurrently busy to the string. If the device is capable of DLSE mode, four channel paths can be concurrently busy to the string.

CPExpert use an M/M/C queuing model (Erlang's C formula) to compute the probability that all paths to the device were busy when the device attempted to reconnect. The M/M/C queuing model is adjusted depending upon the number of concurrent paths¹. The average number of missed RPS reconnect attempts is given by the formula:

$$N = \frac{U}{1 - U}$$

where U is the probability that a reconnect attempt will find all paths busy

The average time spent attempting to reconnect is simply the rotational time of the device multiplied by the average number of missed RPS reconnect attempts. This yields an estimate of the average I/O delay caused by missed RPS reconnect attempts.

¹ Please refer to *Probability, Statistics, and Queuing Theory* by Arnold O. Allen for a description of the M/M/C queuing model.

Rule DAS120 is produced if the average time spent attempting to reconnect to the channel path accounts for a significant percent of the device response time.

The I/O delay caused by missed RPS reconnect may be under-estimated if the device is shared between systems. This is because the controller may be busy to a different system when the device attempts to reconnect. There is no information in RMF records to show the time the controller was busy to another system. RMF data from all systems sharing the controller must be consolidated and analyzed in order to obtain this information. The DASD Component does not perform this analysis at present. However, the logic necessary to determine this situation will be added in future updates to the product.

Suggestion: Missed RPS reconnect time is caused by too much activity on the channel paths to the device. This problem can be corrected by:

- Removing I/O data transfer from the path(s). When devices are seeking or searching for a sector, the channel path is not busy. Therefore, high channel activity is primarily due to transferring data (controller protocol accounts for some path busy, but this time generally is quite small). Removing I/O data transfer from the path(s) can be accomplished by:
 - Moving data sets from the paths. This often is the easiest solution, since data sets with high data transfers can be relocated to other volumes on different channel paths.
 - The volumes responsible for high data transfer time could be moved to other strings on different channel paths.
- Rescheduling workloads to minimize the contention. For example, some batch jobs may be performing heavy I/O to volumes on the string. These jobs may be rescheduled to a time when their I/O would not cause problems.
- Adding paths to the device. If the device is not capable of DLS or DLSE mode, this will require upgrading the device (e.g., from IBM-3380 Model A04 to a more recent model). If the device is capable of DLS mode (that is, it can dynamically reconnect to either of two paths), you may consider upgrading the device to one capable of DLSE mode (that is, it can dynamically reconnect to any of four paths). This action may require that you upgrade your controller from an IBM-3880 storage controller to an IBM-3990 storage controller. Since upgrades are relatively expensive, you should first assess the feasibility of moving active data sets or volumes from the path.

-
- Adding channel paths, acquiring additional controllers, and moving some volumes to the new controllers. This would be a more expensive option, and may not be feasible (depending upon the processor and I/O configuration).

Rule DAS121: VOLUMES CONTRIBUTING TO RPS DELAY

Finding: CPEXpert identifies the volumes contributing to missed RPS reconnect delays.

Impact: This information can be useful when deciding on a course of action to correct the missed RPS reconnect problems. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).* |

Logic flow: The following rules cause this rule to be invoked:
 DAS100: Volume with the worst overall performance
 DAS120: Missed RPS reconnect was major cause of I/O delay

Discussion: If Rule DAS120 is produced, CPEXpert examines the SMF Type 74 information and IOCP macro information to select all devices sharing paths with the device experiencing missed RPS reconnect delays. Devices contributing 5% or more to the utilization of the channel paths are listed, ordered descendingly by their contribution to path utilization.

A device contributes to path utilization mostly based upon the connect time of the device. Consequently, CPEXpert displays the devices having the largest per-second connect times on the path.

Suggestion: You should consider separating the volumes contributing to missed RPS delays from those on the volume experiencing the significant missed RPS delays.

- You can move data sets from the volumes contributing the most to path utilization.
- You can move volumes to a different string.
- You might reschedule workload to minimize the path utilization at critical times.

Rule DAS123: NON-DASD DEVICES CONTRIBUTED TO RPS DELAY

Finding: CPExpert has determined that non-DASD I/O devices were attached to a channel path of the volume experiencing missed RPS reconnect delays. These non-DASD I/O devices were busy a significant percent of the time and contributed to the RPS delay.

Impact: This finding can have a HIGH IMPACT on the performance of the device experience the missed RPS reconnect delays. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).* |

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance
DAS120: Missed RPS reconnect was major cause of I/O delay

Discussion: CPExpert determines whether any non-DASD I/O devices (e.g., tapes drives, etc.) share channel paths with DASD. If missed RPS reconnect delays were a major cause of I/O delay, CPExpert undertakes an analysis of the non-DASD I/O devices sharing channel paths. CPExpert examines the SMF Type 74 information to determine whether these non-DASD devices had a significant connect time to the path.

CPExpert uses a M/M/c queuing model to estimate the amount of missed RPS reconnect delay caused by the path utilization of the non-DASD devices.

Rule DAS123 is produced if the queuing model estimates that path utilization of the non-DASD devices causes more than 10% of the missed RPS reconnect delay.

Suggestion: CPExpert suggests that you eliminate or minimize the impact of the non-DASD I/O devices on the DASD performance by:

- Consider rescheduling the workload accessing the non-DASD I/O devices to a period when the data transfer would not cause DASD problems.
- Remove the non-DASD I/O devices from the channel paths serving the DASD devices. This may mean that you must acquire additional channel paths.

-
- If neither of the above options are feasible, consider placing only low-utilization (and non-critical) DASD on the paths shared with the non-DASD I/O devices.

Rule DAS125: APPLICATIONS CONTRIBUTING TO RPS DELAY

Finding: CPExpert identifies the applications contributing to missed RPS reconnect delays to the volume with worst performance.

Impact: This information can be useful when deciding on a course of action to correct the missed RPS reconnect problems. This rule applies only if "expanded" analysis is performed (that is, only if the modification is made to MXG or MICS so that volume information can be associated with workloads). *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).*

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance
DAS120: Missed RPS reconnect was major cause of I/O delay

Discussion: If Rule DAS120 is produced, CPExpert examines the SMF Type 30(Interval) information to select all applications that reference the volume. An application contributes to path utilization (and thus, contributes to missed RPS reconnect) mostly based upon the connect time of I/O operations to the device.
CPExpert lists the applications accessing the volume with the worst performance, during the entire measurement period being analyzed. The list is ordered descendingly by the average percent use of paths, so you can identify the applications which probably had the most impact on the volume.

As described in Section 5, the SMF Type 30(Interval) information is not synchronized with the SMF Type 70(series) information. Consequently, the identification of applications may not correctly identify the applications with the most I/O operations to the volume.

However, if the analysis produces the same results after analyzing more than one day's measurement data, you can be more comfortable that the applications are correctly identified.

Suggestion: In many cases, once the applications are identified, either applications personnel or systems personnel will verify whether the particular applications are responsible for the majority of the I/O operations to the particular volume.

You should become comfortable that the applications presented are responsible for the utilization of the paths during the period when missed RPS reconnect was a major cause of device delay. You can then take action to (1) reschedule the applications to minimize contention, (2) examine the files referenced by the applications, or (3) change the way in which the applications access their files.

Rule DAS130: PEND TIME WAS MAJOR CAUSE OF I/O DELAY

Finding: CPExpert has determined that PEND time was a major cause of delay in DASD response for the device.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance

Discussion: PEND time is the time from the issuance of the StartSubChannel (SSCH) instruction until the device is selected by the control unit and physical positioning commands (such as seek and set sector, or define extent) are transferred to the device.

With modern fixed block architecture (FBA) devices, the PEND time ends when the physical positioning commands are presented to the *logical volume control block* within the control unit. The PEND time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, or wait for device, or wait for “other” reasons)¹.

PEND is measured by the channel subsystem. After IOS issues the Start Subchannel command, the channel subsystem may not be able to initiate the I/O operation if any path or device busy condition is encountered:

C The channel selected for the I/O operation could be busy with another I/O operation from another system image in the same CEC. This time is not reflected in the SMF data.

C The director port could be busy with another I/O operation². This time is reflected in SMF data as SMF74DPB.

C The control unit could be busy with another I/O operation from another system. This time is reflected in SMF data as SMF74CUB.

¹PEND time is significantly reduced with FICON channels. FICON channels can have multiple I/O operations concurrently active, which reduces the potential PEND time caused by channel busy. There is no port busy time with FICON switches, and control unit time is significantly reduced. This statement regarding PEND time is not necessarily correct if a large number (more than 5) I/O operations are concurrently executing on a FICON channel. Dr. H. Pat Artis and Mr. Robert Ross have presented the results of research indicating that performance can degrade significantly when more than 5 I/O operations (Open Exchanges) are concurrently active on a FICON channel (see “Understanding FICON Channel Path Metrics” at www.perfassoc.com).

²Director port busy can occur only on an ESCON channel. The use of buffer credits on a FICON native channel eliminates director port busy.

C The device could be busy with I/O from another system. This time is reflected in SMF data as SMF74DVB.

There can be “other” PEND time not reflected in the above descriptions. For many systems, “other” PEND time is zero or very small. For some systems, the “other” PEND time is dramatically large (often, 75% or more of the average response time).

One possible cause of the “other” PEND time is PEND for channel busy. If all channels between the MVS image and the device are busy, the channel subsystem must wait until a channel becomes available. This wait for channel is reflected in PEND time. Depending on the number of MVS images using the channels to the device, channel activity could be high. This activity could (and often would) be caused by activity to other logical volumes, rather than the device exhibiting poor performance.

As mentioned earlier, PEND for channel busy is not reflected in the SMF data³. However, CPEXpert calculates an estimated PEND for channel busy based on I/O configuration information.

When CPEXpert creates the model of the I/O configuration, it retains information about each path to a device. Included in this path information is the physical path busy at the CEC level, for each path. Consequently, CPEXpert has an overall view of all physical paths to the device, and can calculate overall channel activity for all channels to the device.

After computing an estimated PEND for channel busy, CPEXpert computes “other” PEND by subtracting estimated PEND for channel, PEND for director port, PEND for control unit, and PEND for device from the total Device Pending time contained in SMF Type 74 records for a particular device.

At present, there is only conjecture⁴ about additional cause of this “other” PEND time. Perhaps either IBM will better describe this “other” PEND time in future, or perhaps research will reveal likely causes of the “other” PEND time.

PEND time can be significant with shared systems. If one system does an I/O request to a device while the storage subsystem is already processing an I/O to that device that came from another system, then the storage

³MXG contains a variable AVGPNCHA (titled “AVG (MS)*PEND DUE TO*CHANNEL BUSY”). However, the MXG AVGPNCHA variable is simply created from the AVGPNDIR (“AVG (MS)*PEND DUE TO*DIRECTOR PORT”) variable. MICS does not contain a “PEND CAUSED BY CHANNEL BUSY” variable.

⁴According to MXG (ADOC74 comments), Dr. H. Pat Artis believes that the “other” PEND is often the internal response time of the subsystem, i.e., the time it takes the subsystem to accept, validate, and acknowledge the first Channel Control Word (CCS) of the channel program.

subsystem will send back a *device busy* indication, resulting in PEND time. This delays the new request and adds to processor and channel overhead.

CPEXpert computes the average per-second PEND delay time, for each of the causes listed above. Rule DAS130 is produced if the average PEND time accounted for a significant percent of the device response time.

The following example illustrates the output from Rule DAS130:

RULE DAS130: MAJOR CAUSE OF I/O DELAY WAS PEND TIME.

A major cause of the I/O delay with VOLSER PPVOL1 was PEND time. The average per-second PEND delay for I/O is shown below:

MEASUREMENT INTERVAL	PEND CHAN	PEND DIR PORT	PEND CONTROL	PEND DEVICE	PEND OTHER	TOTAL PEND
8:30- 8:45,22OCT2001	0.013	0.000	0.000	0.003	0.142	0.158
8:45- 9:00,22OCT2001	0.018	0.000	0.000	0.004	0.205	0.226

Suggestion: No suggestions are associated with this finding. CPEXpert will analyze the high PEND delay time. The following rules will be produced to indicate which is the major cause:

C DAS131: PEND time was caused by channel busy

C DAS132: PEND time was caused by director port busy

C DAS133: PEND time was caused by controller busy delays

C DAS134: PEND time was caused by device busy delays

C DAS135: PEND time was caused by other delays

Please note that the DAS132-DAS135 rules are "LEVEL-2" rules, which means that they will not be produced unless you have specified %LET VERBOSE=VERBOSE in USOURCE(DASGUIDE). If the rules are not produced, you can simply examine data associated with Rule DAS130 to select the major cause. Then you can examine the documentation associated with the major cause.

Rule DAS131 is a "LEVEL-1" rule because it provides additional information on physical channel busy times, by CHPID.

Rule DAS131: PEND DELAY TIME WAS CAUSED BY CHANNEL ACTIVITY

Finding: CPExpert has determined that the excessive PEND time was caused by utilization of the channels to the device.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance
DAS130: Major cause of I/O delay was PEND time

Discussion: PEND time is the time from the issuance of the SSCH instruction until the device is selected by the control unit. This time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, wait for device, or wait for other reasons).

With modern fixed block architecture (FBA) devices, the PEND time ends when the physical positioning commands are presented to the *logical volume control block* within the control unit. The PEND time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, or wait for device, or wait for “other” reasons)¹.

The PEND time due to channel busy means the channel that was selected for the I/O operation was busy with another I/O operation from another system image in the same CEC. Since the channel was busy, the SSCH instruction could not result in device selection, and a PEND condition existed.

SMF Type 74 records do not contain the PEND time caused by channel busy². However, CPExpert calculates an estimated PEND for channel busy based on I/O configuration information.

¹PEND time is significantly reduced with FICON channels. FICON channels can have multiple I/O operations concurrently active, which reduces the potential PEND time caused by channel busy. There is no port busy time with FICON switches, and control unit time is significantly reduced. This statement regarding PEND time is not necessarily correct if a large number (more than 5) I/O operations are concurrently executing on a FICON channel. Dr. H. Pat Artis and Mr. Robert Ross have presented the results of research indicating that performance can degrade significantly when more than 5 I/O operations (Open Exchanges) are concurrently active on a FICON channel (see “Understanding FICON Channel Path Metrics” at www.perfassoc.com).

²MXG contains a variable AVGPNCHA, which is titled 'AVG (MS)*PEND DUE TO*CHANNEL BUSY'. However, the AVGPNCHA variable is simply created from the AVGPNDIR titled 'AVG (MS)*PEND DUE TO*DIRECTOR PORT' variable. MICS does not contain a “PEND CAUSED BY CHANNEL BUSY” variable.

When CPExpert creates the model of the I/O configuration, it retains information about each path to a device. Included in this path information is the physical path busy at the CEC level, for each path. Consequently, CPExpert has an overall view of all physical paths to the device, and can calculate overall channel activity for all channels to the device.

Rule DAS131 is produced when the calculated PEND time due to channel busy accounts for more than one-third of the device PEND time. This output will be produced for all channels sets (by CEC serial number) that are used to reference the logical volume experiencing high PEND time.

The following example illustrates the output from Rule DAS131:

```

RULE DAS131:  LARGE PEND TIME DELAY WAS CAUSED BY CHANNEL BUSY.

A significant amount of the PEND time delay was caused by high channel
utilization for the channels connected to VOLSER CICS11.  This volume
was referenced by the indicated channels.

                                AVERAGE PHYSICAL CHANNEL BUSY FOR CHPIDS:
MEASUREMENT INTERVAL          3D  59  67  72  7F  99  A7  B4
8:30- 8:45,22OCT2001          19  18  84  33  53  42  31  23

```

Suggestion: If an important device is experiencing delays because of high PEND caused by channel utilization, you should consider the following alternatives:

C **Redistribute data sets.** The high PEND time might be solved by redistributing high activity data sets among different volumes on different paths.

If SMF Type 42 (Data Set Statistics) are available, CPExpert will identify data sets on the logical volume that have heavy I/O activity. However, please keep in mind that the PEND time is caused by channel activity. The I/O activity of the particular volume experiencing high PEND time might not be (and probably is not) the cause of high PEND time. Consequently, examining the results for the Type 42 (Data Set Statistics) for the volume might not yield satisfactory results.

C Move the logical volume to a different controller referenced by different channels.

C If redistributing the data sets or moving the volume is not feasible, perhaps more channels can be assigned to the logical control unit through which the device is referenced.

Rule DAS132: PEND DELAY TIME WAS CAUSED BY DIRECTOR PORT

Finding: A significant amount of the PEND time delay was caused by director port busy.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
 DAS100: Volume with the worst overall performance
 DAS130: Major cause of I/O delay was PEND time

Discussion: PEND time is the time from the issuance of the SSCH instruction until the device is selected by the control unit. This time is caused by queuing for the path (wait for channel, wait for control unit or wait for head-of-string), and can be caused by other systems sharing the device (wait for device).

The PEND time due to director port busy, means the director port was busy with another I/O operation from another system image in the same CEC. Since the director port was busy, the SSCH instruction could not result in device selection, and a PEND condition existed.

SMF Type 74 records contain the PEND time caused by director port busy (variable SMF74DPB). This variable is contained in MXG as AVGPNDIR and in MICS as DBADPBTM.

CPEXpert produces Rule DAS130 to report the causes of PEND delay time. Rule DAS132 is produced when director port busy delay was the major cause of PEND delay time.

Suggestion: If Rule DAS132 is consistently produced, you should consider the following alternatives:

C Redistribute data sets. The high PEND time might be solved by redistributing high activity data sets among different volumes on different paths.

If SMF Type 42 (Data Set Statistics) are available, CPEXpert will identify data sets on the logical volume that have heavy I/O activity. However, please keep in mind that the PEND time is caused by director port activity. The I/O activity of the particular volume experiencing high PEND

time might not be (and probably is not) the cause of high PEND time. Consequently, examining the results for the Type 42 (Data Set Statistics) for the volume might not yield satisfactory results.

C Move the logical volume to a different controller referenced by different channels.

C If redistributing the data sets or moving the volume is not feasible, then perhaps more channels can be assigned to the logical control unit through which the device is referenced.

Rule DAS133: PEND DELAY TIME WAS CAUSED BY CONTROLLER BUSY

Finding: A significant amount of the PEND time delay was caused by controller activity.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
 DAS100: Volume with the worst overall performance
 DAS130: Major cause of I/O delay was PEND time

Discussion: PEND time is the time from the issuance of the SSCH instruction until the device is selected by the control unit. This time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, wait for device, or wait for other reasons).

Large PEND times for devices attached to cached controllers may imply a high percent of read miss operations, or non-volatile storage (NVS) writes.

To improve the probability of a read hit, the controller can *prestage* data into its cache. Prestaging means that data is read into the controller's cache ahead of its actually being required for use by an application. The amount of data that is prestaged depends on (1) whether the data is being accessed in a direct (random) mode or in a sequential mode and (2) the controller model and the enhancements made to the controller.

C For *direct mode*, after the record is located, the 3390-3 and 3990-6 (initial version) stages in the balance of the track being read.

The 3990 Model 6 (with record cache) stages only the records requested into cache, eliminating the balance of the track staging that is normal with track caching as was implemented on initial versions of 3990-6 and on the 3990-3. This improvement reduces the PEND time caused by the controller busy during track staging.

C As examples of prestaging for *sequential mode*, the 3990-3 reads up to two tracks into the cache¹ before they are required, while the ESS 2105 sequential staging reads up to two cylinders ahead.

¹With the Sequential Staging Performance Enhancement, the 3990-3 can prestage up to a full cylinder (15 tracks) into the cache.

During prestaging operations for sequential reads, the control unit regularly checks to see whether other I/O requests are waiting to be processed. If any are waiting, the control unit interrupts the prestage operation, processes the queued requests, and continues with the prestage.

In DASD Fast Write Mode, the data is stored simultaneously in cache storage and in nonvolatile storage (NVS). At some subsequent time, the data in NVS can be *destaged* to DASD.

In Cache Fast Write Mode, data is placed into cache immediately, and there is no interaction with the device nor with NVS. However, if cache memory is required (or if Cache Fast Write Mode is turned off), the data in cache is destaged to DASD.

Significant PEND time can result from destaging to DASD.

Rule DAS134: PEND DELAY TIME WAS CAUSED BY DEVICE BUSY

Finding: A significant amount of the PEND time delay was caused by device busy.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
 DAS100: Volume with the worst overall performance
 DAS130: Major cause of I/O delay was PEND time

Discussion: PEND time is the time from the issuance of the SSCH instruction until the device is selected by the control unit. This time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, wait for device, or wait for other reasons). Some of the causes of PEND for device busy are listed below:

C PEND time for device busy can be caused by other systems in the sysplex that issue a RESERVE for the device. While the RESERVE is held, I/O operations to the device will be held in a PEND for device busy state.

Additional information will be available if CPExpert is performing shared DASD analysis (shared DASD analysis will be done if %LET SHARED=Y; was specified in USOURCE(DASGUIDE) and if the SMF data indicates that the device was shared). Rule DAS300 will be produced to show the effects of activity from other systems on the device being analyzed. Included in the output from Rule DAS300 will be the amount of RESERVE time from each system, and the total from all systems sharing the device.

C Multiple Allegiance allows multiple active concurrent I/O operations on a particular device when the I/O requests originate from different systems. With Multiple Allegiance, there is complete access with read I/O operations. For write I/O operations, there is concurrent access unless there is a conflicting extent¹. If there is a conflicting extent, the controller holds the I/O operation in a PEND state for the device.

C After an I/O operation, the device will read the remainder of the track into its device-level buffer. This is done to prevent delay for rotational positioning. If a new I/O operation is attempted while data is being read

¹ A conflicting extent is one in which the write operation attempts to update an extent.

into the device cache buffer, the I/O operation will be in a PEND state for device busy.

Suggestion: If Rule DAS134 is produced frequently, you should consider the following alternatives.

C If shared device analysis is specified as **%LET SHARED = Y;** in USOURCE(DASGUIDE), CPExpert will analyze potential problems caused by sharing DASD. Rule DAS300 will be produced for all systems that share the device with high PEND time, if CPExpert concludes that other systems could cause performance problems. The RESERVE time will be included in the output from Rule DAS300.

If this RESERVE time is high for the device, you should consider whether high activity data sets can be moved among different volumes on different paths.

C Alternatively, determine whether the data sets can be moved to a controller that supports Multiple Allegiance². Multiple Allegiance allows multiple active concurrent I/O operations on a particular device when the I/O requests originate from different systems. With Multiple Allegiance, there is complete access with read I/O operations. For write I/O operations, there is concurrent access unless there is a conflicting extent.

C Alternatively, consider whether workload scheduling can eliminate the conflicts between the data access requirements between systems.

²Multiple Allegiance is available with IBM's Enterprise Storage Server (ESS) subsystems.

Rule DAS135: PEND DELAY TIME WAS CAUSED BY OTHER DELAYS

Finding: A significant amount of the PEND time delay was caused by delays that were not unexplained by either the causes of PEND time as reported by SMF, or estimated causes as calculated by CPEXpert.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
 DAS100: Volume with the worst overall performance
 DAS130: Major cause of I/O delay was PEND time

Discussion: PEND time is the time from the issuance of the SSCH instruction until the device is selected by the control unit. This time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, wait for device, or wait for other reasons).

SMF contains information describing the PEND delay time caused by wait for director port (SMF74DPB), wait for control unit (SMF74CUB), and wait for device (SMF74DVB). As described in Rule DAS131, CPEXpert computes an estimated PEND time caused by channel busy.

After computing an estimated PEND for channel busy, CPEXpert subtracts the wait for channel, wait for director port, wait for control unit, and wait for device from the total Device PEND time reported by SMF in SMF74PEN. After this subtraction, there often is a remainder. This remainder¹ has been titled "other" PEND time since it is PEND time that is not reflected in the measured or estimated causes of PEND delay time.

For many systems, "other" PEND time is zero or very small. For some systems, the "other" PEND time is dramatically large (often, 75% or more of the average response time).

Suggestion: There are no suggestions with Rule DAS135. The finding is produced simply to alert you that the PEND delay data contains a large amount of

¹MXG does not calculate the PEND due to channel busy, but simply sets the AVGPNDCHA equal to the AVGPNDIR. MXG does subtract the PEND for director port, PEND for control unit, and PEND for device from the total device PEND delay. Since MXG does not calculate an estimated PEND time due to channel busy, MXG produces a much larger PEND "other" value than does CPEXpert's calculations.

delay that was unexplained by either the causes of PEND time as reported by SMF, or estimated causes as calculated by CPExpert.

At present, there is only conjecture² about additional cause of this “other” PEND time. Perhaps either IBM will better describe this “other” PEND time in future, or perhaps research will reveal likely causes of the “other” PEND time.

²According to MXG (ADOC74 comments), Dr. H. Pat Artis believes that the “other” PEND is often the internal response time of the subsystem, i.e., the time it takes the subsystem to accept, validate, and acknowledge the first Channel Control Word (CCS) of the channel program.

Rule DAS140: CONNECT TIME WAS A MAJOR CAUSE OF I/O DELAY

Finding: Connect time was a major cause of the I/O delay with the volume.

Impact: This finding may have a LOW IMPACT or MEDIUM IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
 DAS100: Volume with the worst overall performance

Discussion: Connect time is the time in which the device is actually connected to the path. This time includes the data transfer time, but also includes protocol exchange (or "hand shaking") between the various components at several stages of the I/O operation.

The data transfer time obviously is a function of the amount of data being transferred. This simply is the number of bytes transferred divided by the transfer speed (for example, if 4096 bytes were transferred from an IBM-3380 with a transfer speed of 3,000,000 bytes per second, the 4096 bytes would require $4096/3,000,000$ seconds; or about 1.36 milliseconds).

Large connect times generally are caused by the following situations:

- A large average block size. This situation may be highly desirable for sequential data sets, but would be undesirable for randomly accessed data.
- Long multi-track searches. For example, the catalog must be searched for cataloged files, the Volume Table of Contents (VTOC) must be searched to find a requested file, a directory must be searched for partitioned data sets, etc.. These searches will result in long connect times for the volume involved.
- Program loading from system packs.

Rule DAS150: QUEUING IN IOS WAS A MAJOR CAUSE OF I/O DELAY

Finding: Queuing in the I/O Supervisor (IOSQ) was a major cause of the I/O delay with the volume.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance

Discussion: IOSQ time is the time from the issuance of a STARTIO macro until the Start SubChannel (SSCH) instruction is issued. After the STARTIO macro is issued, the software determines whether the device is busy with the system on which the STARTIO macro was issued (that is, whether there is an available Unit Control Block (UCB) for the device). If the device is not busy with this system, the SSCH instruction is issued. However, if the device is busy with this system (a UCB is available), the I/O request is queued. Thus, IOSQ time always means that the device is unable to handle additional requests from this system.

This discussion of IOSQ time does not always apply to Parallel Access Volumes (PAVs)¹. With PAV devices, MVS creates multiple UCBs for each device, depending on how many “alias devices” have been defined. The multiple UCBs allow multiple active concurrent I/Os on a given device when the I/O requests originate from the same system². Using PAVs can dramatically improve I/O performance by nearly eliminating IOSQ.

Rule DAS150 is produced if the average IOSQ time accounted for a significant percent of the device response time.

The following example illustrates the output from Rule DAS150:

¹PAV devices are available with Enterprise Storage Server (ESS). With PAV devices, a “base device” address is defined, and a UCB is associated with this base address. “Alias device” addresses can be defined and UCBs are associated with the alias device addresses.

²Multiple Allegiance allows multiple active concurrent I/O operations on a given device when the I/O requests originate from different systems. The Multiple Allegiance feature is available with Enterprise Storage Server (ESS).

RULE DAS150: MAJOR CAUSE OF I/O DELAY WAS QUEUING IN I/O SUPERVISOR.

A major cause of the I/O delay with VOLSER SP0006 was queuing in the I/O Supervisor (IOS). Please refer to the DASD Component User Manual for a discussion of ways to reduce I/O queuing;

MEASUREMENT INTERVAL	AVERAGE REQS IN QUEUE	AVG IOSQ DELAY PER I/O	PAV BASE DEVICE	PAV UCB COUNT	UCB COUNT CHANGED
8:30- 8:45,22OCT2001	3.8	0.121	N		
8:45- 9:00,22OCT2001	2.2	0.087	N		

Suggestion: Large IOSQ times usually involve the following situations:

- Multiple data sets may be active on the volume. This situation is the most common and easiest to solve. The data sets can be redistributed among different logical volumes, to eliminate the queuing for the single volume.
- The data sets can be placed on PAV devices or redistributed among different logical volumes, to eliminate the queuing for the single volume.
- If using static PAVs, assign more aliases to the device.
- If using dynamic PAV, increase the number of PAVs associated in the pool for the subsystem.
- Ensure that all PAVs that should be bound to the device are online and are operational. You can use the DEVSERV QP and DS QP,xxxx,UNBOX commands to do this.
- Multiple users may be using the same data set on the volume. Depending upon the data set characteristics, duplicate copies of the data set placed on different volumes may solve the IOSQ problems.
- Multiple application systems may be using the volume experiencing high IOSQ times. In this case, perhaps application redesign or scheduling can solve the problem.
- A particular application (or system function) may be executing I/O to the device faster than the device can respond.
- The overall device response time (PEND, DISC, and CONN) times may be large, such that the device is unable to provide quick response to the I/O requests. This situation will be revealed by large values in the PEND, DISC, or CONN measures. Consider moving files to a faster storage (coupling facility structure, expanded storage, Data In Memory, etc.). Also, consider speeding up or reducing the I/O on the path or the device (e.g., specify optimal VSAM options, revise blocking options, etc.).

Rule DAS151: QUEUING WAS PROBABLY CAUSED BY INDEPENDENT APPLICATIONS

Finding: More than half of the I/O queuing to the device was explained using a queuing model. Consequently, CPEXpert believes that the queuing probably was caused by independent applications.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device. If the finding is correct, and the independent applications are referencing different files, then the corrective actions could result in significant performance improvements.

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance
DAS150: Queuing was the major cause of response delay

Discussion: CPEXpert uses a M/M/1 queuing model to calculate an estimated queue time for each measurement interval being analyzed. The underlying assumptions of the model are exponential interarrival times, exponential service distributions, and an infinite population. If device activity occurs in this way, the queuing model can predict the expected queuing delays.

If the queuing delay as measured by RMF is significantly different from the estimated delay from the model, it would be clear that the activity did not occur in a random fashion, and most likely the cause of the difference would be that the interarrival times are not randomly distributed.

However, the I/O delay to this pack was fairly well explained by the queuing model. More than half of the I/O queuing delay was explained by the model, for a majority of the measurement intervals. This indicates that the activity was mostly random, as would be expected from independent applications accessing files on the pack.

Suggestion: The most improvement for this pack would likely result from (1) separating the files to different packs, (2) rearranging the files within the pack, (3) tuning the file structure (for example, compressing a shared partitioned data set (PDS), (4) scheduling the applications to avoid contention, or (5) examining the applications doing most of the I/O.

If CPEXpert is performing "expanded" analysis, Rule DAS180 will provide information regarding the applications referencing the volume. The results

from this rule should be examined to determine which applications are candidates for actions.

You can identify the files with the most activity in one of several ways, depending upon the volume and the application involved. The method you select should depend upon the availability of information and tools in your environment.

- Discuss the file activity with applications personnel or systems personnel (depending upon the volume). If knowledgeable personnel are available, this is often the quickest and easiest way to determine how to solve the problems. Once a file access problem has been brought to their attention, applications or systems personnel often can easily decide upon a good solution.
- Use an exit available in MXG or in MICS to select Type 30(DD) information just for the volume. The Type 30(DD) information is rarely retained in a performance data base because it is so voluminous. However, you easily can code an exit in MXG or MICS to select the Type 30(DD) information for a specific volume. The amount of data selected would not be too large in most cases.

You can then write a SAS program to list the DD names used by applications, weighted by the number of accesses. For example, you could code the following to analyze data extracted during MXG update processing:

```
PROC FREQ DATA=pdblib.file;  
  WHERE DEVNR = addressX;  
  TABLES DDNAME/OUT=TEMP NOPRINT;  
  WEIGHT EXCPS;  
PROC SORT DATA=TEMP;  
  BY DESCENDING COUNT;  
PROC PRINT;  
RUN;
```

The device address is displayed by RULE DAS100, so the "address" in the above coding would be replaced with the actual address displayed by that rule. If you have a MXG performance data base, the address is retained in hexadecimal format, so you would suffix an X to the address. If you have a MICS performance data base, the variable names must be changed appropriately (for example, DEVNR would be changed to DEVADDR). Additionally, the address is retained in MICS as a character representation of the value, so you do not need to suffix an X to the address. Rather, you must enclose the address in quotes.

The result from the above code would be a list of all DDNAMEs referencing the device being analyzed, weighted by the number of EXCPs to the device, and ordered descendingly with the DD statements for the most active files listed first. You often cannot be certain that the most referenced files are causing problems. However, in most cases, you will find that only a few files (often only 2 or 3) account for over 90% of the accesses. These files generally will be the ones causing arm contention problems.

You can then write a SAS program to list the applications (jobs) referencing the device, weighted by the number of accesses. For example, you could code the following to analyze data extracted during MXG update processing:

```
PROC FREQ DATA=pdblib.file;  
  WHERE DEVNR = addressX;  
  TABLES JOB/OUT=TEMP NOPRINT;  
  WEIGHT EXCPs;  
PROC SORT DATA=TEMP;  
  BY DESCENDING COUNT;  
PROC PRINT;  
RUN;
```

The result from the above code would be a list of all jobs referencing the device being analyzed, weighted by the number of EXCPs to the device, and ordered descendingly with the jobs with the most active files on the device listed first. These jobs generally will be the ones causing arm contention problems. (If you have a MICS performance data base, you would change the code as described earlier.)

- Use a commercially-available DASD activity monitor to isolate the files accessed on the problem volume. If such a monitor is available, this would be a direct way to determine the problems.

Rule DAS152: QUEUING WAS PROBABLY CAUSED BY A SINGLE APPLICATION

Finding: Less than half of the I/O queuing to the device was explained using a queuing model. Consequently, CPEXpert believes that the seeking probably was caused by a single application, rather than by independent applications.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device. If the finding is correct, then actions directed to the specific application could result in significant performance improvements.

Logic flow: The following rules cause this rule to be invoked:
 DAS100: Volume with the worst overall performance
 DAS150: Queuing was the major cause of response delay

Discussion: CPEXpert uses a M/M/1 queuing model to calculate an estimated queue time for each measurement interval being analyzed. The underlying assumptions of the model are exponential interarrival times, exponential service distributions, and an infinite population. If device activity occurs in this way, the queuing model can predict the expected queuing delays.

If the queuing delay as measured by RMF is significantly different from the estimated delay from the model, it would be clear that the activity did not occur in a random fashion, and most likely the cause of the difference would be that the interarrival times are not randomly distributed.

This was the case with the volume being analyzed. Consequently, CPEXpert concludes that the activity was caused by a single application accessing the device. Note that there could be more than one application accessing the device, but if their access patterns were not random with respect to each other, then this conclusion would still be valid.

Suggestion: The most improvement for this pack would likely result from (1) separating the files to different packs, (2) rearranging the files within the pack, (3) tuning the file structure (for example, compressing a shared partitioned data set (PDS), or (4) examining the application doing most of the I/O.

You can identify the files with the most activity in one of several ways, depending upon the volume and the application involved. The method you

select should depend upon the availability of information and tools in your environment.

- Discuss the file activity with applications personnel or systems personnel (depending upon the volume). If knowledgeable personnel are available, this is often the quickest and easiest way to determine how to solve the problems. Once a file access problem has been brought to their attention, applications or systems personnel often can easily decide upon a good solution.
- Use an exit available in MXG or in MICS to select Type 30(DD) information just for the volume. The Type 30(DD) information is rarely retained in a performance data base because it is so voluminous. However, you easily can code an exit in MXG or MICS to select the Type 30(DD) information for a specific volume. The amount of data selected would not be too large in most cases.
- If CPExpert is performing "expanded" analysis, Rule DAS180 will provide information regarding the applications referencing the volume. The results from this rule should be examined to determine which applications are candidates for actions.

Rule DAS153: WORST QUEUING WAS PROBABLY CAUSED BY INDEPENDENT APPLICATIONS

Finding: CPEXpert determined that seeking was the major cause of delay in DASD response for the device during the measurement interval with the worst I/O response. More than half of the I/O queuing to the device was explained using a queuing model. Consequently, CPEXpert believes that the queuing during this interval probably was caused by independent applications.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device. If the finding is correct, and the independent applications are referencing different files, then the corrective actions could result in significant performance improvements.

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance
DAS150: Seeking was the major cause of response delay

Discussion: CPEXpert uses a M/M/1 queuing model to calculate an estimated queue time for the measurement interval with the worst I/O response for the device being analyzed. The underlying assumptions of the model are exponential interarrival times, exponential service distributions, and an infinite population. If device activity occurs in this way, the queuing model can predict the expected queuing delays.

If the queuing delay as measured by RMF is significantly different from the estimated delay from the model, it would be clear that the activity did not occur in a random fashion, and most likely the cause of the difference would be that the interarrival times are not randomly distributed.

However, the I/O delay to this pack was fairly well explained by the queuing model. More than half of the I/O queuing delay was explained by the model, for a majority of the measurement intervals. This indicates that the activity was mostly random, as would be expected from independent applications accessing files on the pack.

Suggestion: The most improvement for this pack would likely result from (1) separating the files to different packs, (2) rearranging the files within the pack, (3) tuning the file structure (for example, compressing a shared partitioned data set (PDS), (4) scheduling the applications to avoid contention, or (5) examining the applications doing most of the I/O.

The actions to be taken when this rule is produced are the same associated with RULE DAS151. Please refer to that rule for further suggestions.

Rule DAS154: WORST QUEUING WAS PROBABLY CAUSED BY A SINGLE APPLICATION

Finding: CPExpert determined that IOSQ was the major cause of delay in DASD response for the device, during the measurement interval with the worst I/O response. Less than half of the I/O queuing to the device was explained using a queuing model. Consequently, CPExpert believes that the queuing probably was caused by a single application, rather than by independent applications.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device. If the finding is correct, then actions directed to the specific application could result in significant performance improvements.

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance
DAS150: Seeking was the major cause of response delay

Discussion: CPExpert uses a M/M/1 queuing model to calculate an estimated queue time for the measurement interval with the worst I/O response for the device being analyzed. The underlying assumptions of the model are exponential interarrival times, exponential service distributions, and an infinite population. If device activity occurs in this way, the queuing model can predict the expected queuing delays.

If the queuing delay as measured by RMF is significantly different from the estimated delay from the model, it would be clear that the activity did not occur in a random fashion, and most likely the cause of the difference would be that the interarrival times are not randomly distributed.

This was the case with the volume being analyzed. Consequently, CPExpert concludes that the activity was caused by a single application accessing the device. Note that there could be more than one application accessing the device, but if their access patterns were not random with respect to each other, then this conclusion would still be valid.

Suggestion: The most improvement for this pack would likely result from (1) separating the files to different packs, (2) rearranging the files within the pack, (3) tuning the file structure (for example, compressing a shared partitioned data set (PDS), or (4) examining the application doing most of the I/O.

The actions to be taken when this rule is produced are the same associated with RULE DAS152. Please refer to that rule for further suggestions.

Rule DAS160: DEVICE DISCONNECT WAS MAJOR CAUSE OF I/O DELAY

Finding: CPEXpert determined that device disconnect (DISC) time was the major cause of delay in DASD response for the device.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
 DAS100: Volume with the worst overall performance

Discussion: DISC means that there is some delay that is often (but not always) associated with a mechanical movement during which the device disconnects from the control unit (or the control unit disconnects from the channel).

With legacy systems (e.g., 3380 drives attached to 3990-2 control units), the DISC time of most concern was associated with seek (arm movement) and rotational position sensing (time waiting for the disk platter to rotate to the location where desired data resides). Considerable performance improvement efforts were directed at reducing the seek activity and reducing the rotational position sensing (RPS)¹ delays for the legacy systems. These two mechanical delays still exist for most modern *redundant array of independent disks* (RAID)² systems, but their impact can not be directly reduced with normal methods.

With modern disks, data is cached into device cache buffers that contain data read from a track on the disk platter. Using device cache buffers containing the track data eliminated the multiple-RPS delays caused by a path busy when the device tried to reconnect. Required data is read into the device cache buffer during a single rotation and stored until a path is available to transfer the data.

In addition to the cache buffer design, modern control units such as the 3990-6 or 2105 have very large cache memory installed. With cache in the control units, data to be read can be transferred in a variety of ways, depending on where the data resides.

¹RPS delays were caused by a path not being available when the required data came under a device read head. Since a path was not available, the device could not reconnect to the channel or control unit. Consequently, data could not be read and transmitted, and another rotation of the platter was experienced until the data again came under the device read head. Multiple rotations might be required, depending on the busy level of the path.

²An array is an ordered collection of physical devices (disk drive modules) that are used to define logical volumes or devices.

For a read operation, desired data often is found in the control unit cache. If the required data is in cache, the data can be transferred between the control unit cache and the channel, and this transfer is done at channel speed. If the required data is not in cache, the data can be transferred between the device and channel (and concurrently placed into the control unit cache for subsequent access).

For write operations, data can be placed into Non-volatile Storage (NVS) as a part of the control unit. Write operations normally end as the data to be written is placed in the NVS; and the storage processor writes the data to the device asynchronous with other activity (as a “back end” staging operation). See subsequent discussion for more detail about read and write operations.

The storage director can simultaneously transfer data between the channel and device and manage the data transfer of different tracks between the cache and channel, and the cache and the device. With large amounts of cache memory, a high percent of data accesses normally will be resolved from the fast cache memory and the relatively slow device will not cause significant delays.

As a result of the above improvements, DISC time for modern systems is a result of *cache read miss* for read operations, back-end staging delay for write operations, peer-to-peer remote copy (PPRC) operations, and other miscellaneous reasons³. DISC time often can be very small with adequate cache. For example, there would be zero disconnect time for a cache read hit (the record was found in the cache). However, DISC time can be large and can cause serious delay to I/O operations.

C Read operations. With devices attached to cached controllers, a read operation finds required data in the cache (a “read hit”) or does not find required data in the cache (a “read miss”).

If a read operation *finds data in the cache*, acquiring the data involves only the transfer of data from cache. In this case, the data transfer takes place at channel speeds. Channel speeds can vary, depending on the channel type, from about 4.5 MB per second (parallel channels), up to 18 MB per second (ESCON channels), to over 100MB per second (FICON channels).

³ Artis has described a “sibling PEND” condition that results from collisions within the physical disk subsystem of RAID devices. See “Sibling PEND: Like a Wheel within a Wheel,” www.cmg.org/cmgapap/int449.pdf. While this condition is titled “sibling PEND,” the time properly belongs in DISC time, rather than PEND time .

If a read operation *does not find data in the cache*, the data must be read from the physical disk device⁴. With the IBM-3390-3 controller and the initial release of the IBM-3390-6 controller, an entire track would be read into cache for a direct read. This algorithm was changed to read only the record required in a direct read; the change eliminated unnecessary activity by the controller⁵.

The implications of reading the data from the physical disk device differ depending on the type of channel:

- C With parallel channels and ESCON channels, the control unit *disconnects* from the channel while the data is being read. After the data has been read, the control unit attempts to reconnect to the channel. The channel must be available when the control unit attempts to reconnect, or additional overhead results. Consequently, channel busy is an important metric with parallel channels and ESCON channels. IBM suggests that these channel types should not have a consistent busy greater than 50% to avoid unacceptable overhead.
- C With FICON Native channels and control units, the control unit does *not* disconnect from the channel while the data is being read, as disconnect and reconnect protocols have been eliminated with FICON. When the frames of data read from DASD are ready to be presented to the channel, the frames simply queue along with any other frames of data (from other I/O operations transferring data) and the data frames are interleaved at channel transfer rates.

While the device delays caused by cache miss operations do not result in disconnect/reconnect protocol exchanges between channel and control unit, the actual device delay time exists nonetheless⁶. These device delays are timed by a FICON control unit, and the time is reported to RMF as DISC time. Thus, the delay time is available with FICON channels and control units and titled "DISC" time, even though the actual disconnect and reconnect activities do not occur.

⁴The data is read into cache, unless *Inhibit Cache Loading* had been specified. With *Inhibit Cache Loading*, the cache is searched to see whether the record is in cache (from a previous I/O operation). If the requested track is not in cache, the channel program operates directly with DASD. Applications can use *Inhibit Cache Loading* when it is known that records read would not likely be accessed again.

⁵The initial design did not consider that the device and the controller would be "busy" during the transfer of the track from the device to the controller. The belief was that the transfer of the track would be "off line" and not adversely impact performance. However, while the track was being transferred to the controller, the device and controller were busy and other I/O operations were constrained. With very active systems, this constraint could seriously degrade performance. By moving to record-level transfer for direct I/O, this constraint was removed.

⁶This might seem a moot point; if the device delay exists, why should it matter whether the time is a result of disconnect between the channel and control unit or simply device delay time? The difference is that the exchange of disconnect and reconnect protocol traffic between the channel and control unit is eliminated with FICON. This exchange of protocol can add considerable overhead, and it is this overhead that is eliminated with FICON. The FICON controller times the device delays that occur simply for RMF reporting.

In order to improve the probability of a read hit, the controller can *prestage* data into its cache. Prestaging means that data is read into the controller's cache ahead of its actually being required for use by an application. The amount of data that is prestaged depends on (1) whether the data is being accessed in a direct (random) mode or in a sequential mode and (2) the controller model and the enhancements made to the controller.

For *direct mode*, the 3990 Model 6 (with record cache) stages only the records requested into cache, eliminating the balance of the track staging as was implemented on initial versions of 3990-6 and on the 3990-3. As examples of prestaging for *sequential mode*, the 3990-3 reads up to two tracks into the cache⁷ before they are required, while the ESS 2105 sequential staging reads up to two cylinders ahead.

Applications can indicate (using Define Extent) that data is to be processed in a sequential mode. With the 3990-6, IBM included a *sequential detection algorithm* that automatically detects whether data is being read sequentially, even if the user did not indicate that reads were in sequential mode. If the algorithm detects sequential access, data is prestaged automatically. For example, with the ESS 2105, when the algorithm detects that 6 or more tracks have been read in succession, the algorithm triggers the sequential staging process.

During prestaging operations, the control unit regularly checks to see whether other I/O requests are waiting to be processed. If any are waiting, the control unit interrupts the prestage operation, processes the queued requests, and continues with the prestage.

C Write operations. With devices attached to cached controllers, a number of options are available to help improve performance for particular applications. Use of these options vary depending on the data access characteristics of records being written, performance goals associated with the applications, amount of cache and NVS that is available, etc. Some of the common options are Bypass Cache Mode, Normal Caching Mode, Cache Fast Write Mode, and DASD Fast Write Mode.

C Bypass Cache Mode. The Bypass Cache Mode causes the data in the cache to be bypassed. The I/O write request is sent directly to DASD, but a search of the cache is performed because the track in the cache could have been modified by a previous I/O operation. If the track is in the cache, the corresponding cache slot is marked invalid to prevent a read hit by a subsequent I/O operation. If the

⁷ With the Sequential Staging Performance Enhancement, the 3990-3 can prestage up to a full cylinder (15 tracks) into the cache.

cache slot had been modified by a previous cache fast write hit or a DASD fast write hit, the track is destaged and the slot is marked invalid.

The performance of an I/O operation with Bypass Cache Mode is almost the same as if the write were performed via a noncache storage control. The Bypass Cache operation is slightly longer than a write via a noncache controller, because a directory search of the controller's cache is required to determine whether the track is in cache.

The controller presents channel end and device end only after the transfer operation is complete. Since the I/O write operation deals directly with the device, disconnect time can be significant.

The Bypass Cache Mode might be used even though the control unit has considerable cache in situations where low priority files are "cache unfriendly" (meaning that they have a poor locality of reference), with very large files with high write activity when the files might "flood" the cache and cause a low read hit or write hit for other (perhaps more important) file accesses.

- C **Normal Caching Mode.** With Normal Caching Mode, all write I/O commands operate directly with the device. In cache operations without cache fast write or DASD fast write, a write operation follows these general rules⁸:
 - C A format write operates directly with DASD. If the track is in cache, it is invalidated. This ensures that a subsequent read will result in a read miss.
 - C If the track modified by an update write operation is in cache, the cache and DASD are updated concurrently (a write hit). This ensures that the data in cache is current.
 - C If the track modified by an update write operation is not in the cache, the operation is a write miss. Only the data on DASD is updated.
 - C No *new* tracks are transferred from DASD to cache as the result of a write operation.
 - C A track in cache is never made "most recently used" by a write hit in basic caching operations.

⁸Source: IBM's 3990 Planning, Installation, and Storage Administration Guide

If a write hit occurs (the write request updates a record that is already in cache), the controller transfers the data to both cache and DASD. This ensures that the data in cache is current, and is available for a subsequent read operation.

If a write Miss occurs (the write request updates a record that is not in cache) data is transferred from the channel to DASD, and is not placed into cache.

The primary objective of a basic cache write operation is to emulate a DASD write, to ensure that the DASD copy of the data is always valid, and to ensure that any copy of the data retained in cache is valid.

The controller presents channel end and device end only after the transfer operation is complete. Since the I/O write operation deals directly with the device, disconnect time can be significant.

- C **Cache Fast Write Mode.** The Cache Fast Write Mode causes data to be placed into cache immediately, and there is no interaction with the device nor with NVS. Cache fast write is useful in situations where the data that may not be required after the completion of the current job or in situations where the data could be easily reconstructed if necessary (data could be reconstructed if the cache failed).

If the record to be written is already in the cache, this is considered a "write hit" and the entire operation is performed with the cache. With either a write miss (data is not in the cache) or a write hit, no DASD access is required. However, write hits cause the record to be made "most recently used." *When cache space is needed, the controller destages the least recently used data to DASD.*

In most cases when Cache Fast Write Mode is used, the data is only temporary, and can be discarded when no longer required. For example, sorts would not require permanent data for their sort work files.

If the cache is reinitialized, all cache fast write data is lost and the cache fast write identifier is incremented. Subsequent I/O operations with the old cache fast write identifier are reported to the requesting program as a permanent I/O error.

The controller presents channel end and device end after the data has been placed in the cache. Since the I/O write operation deals only

with the cache, disconnect time is eliminated for normal I/O operations⁹.

- C **DASD Fast Write Mode.** In DASD Fast Write Mode, the data is stored *simultaneously* in cache storage and in nonvolatile storage. Since data is stored in NVS, access to a physical DASD is not required for write hits to ensure data integrity. The copy of the data in nonvolatile storage allows storage processor to continue without waiting for the data to be written to DASD. The data remains in cache storage and in nonvolatile storage until the storage control destages the data to DASD. Since completion of the write is indicated when the cache data transfer is complete, DASD Fast Write provides a significant performance enhancement over basic write operations; the DASD fast write hit is as fast as a read hit.

In MVS, activation and deactivation of DASD fast write is provided by a system utilities command with extended function programming support. DASD fast write remains active until explicitly deactivated by another command. DASD fast write is activated at a volume level and is the default for all write operations directed at that volume. DASD fast write can be inhibited at the channel program level.

If DASD fast write is deactivated, the 3990 destages the DASD fast write data to DASD. The 3990 also destages the DASD fast write data to DASD if (1) NVS is deactivated, (2) subsystem caching or device caching is deactivated, and (3) more space is made available in the cache or NVS. These destaging operations are between the cache or NVS and DASD. Consequently, the activity does not result in disconnect time for normal I/O operations (that is, they would not be reflected as DISC time by RMF).

Disconnect delay time is available in SMF Type 74 records. CPExpert produces Rule DAS160 if the average disconnect delay time accounted for a significant percent of the device response time. If RMF Cache Subsystem Activity Statistics records are available, CPExpert produces detailed information about the cache activity for each RMF measurement interval. If RMF Cache Subsystem Activity Statistics are *not* available, CPExpert simply provides the basic message, and an annotation that cache activity data was not available.

The following example illustrates the output from Rule DAS160:

⁹There can be considerable device activity if the data is destaged because cache space was needed or after cache fast write is turned off. This destage activity could adversely impact other I/O operations requiring access to the device.

RULE DAS160: DISCONNECT TIME WAS MAJOR CAUSE OF I/O DELAY.

A major cause of the I/O delay with VOLSER DB0068 was DISCONNECT time. DISC time for modern systems is a result of cache read miss operations, potentially back-end staging delay for cache write operations, peer-to-peer remote copy (PPRC) operations, and other miscellaneous reasons.

MEASUREMENT INTERVAL	----CACHE----		--PERCENT--		DASD	CACHE	PPRC	BPCR	ICLR
	READS	WRITES	READ HITS	WRITE HITS	TO CACHE	TO DASD			
8:30- 8:45,22OCT2001	5998	1286	47.5	99.2	3145	1037	0	1	0
8:45- 9:00,22OCT2001	6215	1375	45.0	99.1	3417	986	0	0	0

The following describes the information presented by Rule DAS160:

- C **Cache reads.** The *cache reads* represent the total read requests. This is the sum of search read caching requests (SMF variable R745DRCR), read sequential requests (SMF variable R745DRSR), and search read cache fast write data requests (SMF variable R745DRNR).
- C **Cache writes.** The *cache writes* represent the total write requests. This is the sum of search read caching requests (SMF variable R745DRCR), read sequential requests (SMF variable R745DRSR), and search read cache fast write data requests (SMF variable R745DRNR).
- C **Percent read hits.** The *percent read hits* represents the sum of search read hits (SMF variable R745DCRH), read sequential hits (SMF variable R745DRSH), and search read cache fast write data hits (SMF variable R745DRNH); divided by the total read requests (described above as “cache reads”).
- C **Percent write hits.** The *percent write hits* represents the sum of search write hits (SMF variable R745DWCH), read sequential hits (SMF variable R745DWSH), and search read cache fast write data hits (SMF variable R745DWNH); divided by the total write requests (described as “cache writes”).
- C **DASD to cache.** The *DASD to cache* represents the normal cache requests (DASD to cache) transfers (SMF variable R745DNTD).
- C **Cache to DASD.** The *cache to DASD* represents the cache to DASD transfers (SMF variable R745DCTD).
- C **PPRC.** The *PPRC* value represents the **Peer to Peer Remote Copy** write count (SMF variable R745XPRC).

C **BPCR**. The BPCR value represents the **Bypass Cache** requests issued (SMF variable R745DBCR).

C **ICLR**. The *ICLR* represents the **Inhibit Cache Loading** requests issued (SMF variable R745DICL).

Rule DAS170: THERE DID NOT APPEAR TO BE A PROBLEM WITH VOLUME PERFORMANCE

Finding: The performance of the volume having the "worst" overall performance was not poor in any area evaluated by CPEXpert.

Impact: This finding has little impact, other than to indicate that there may be no problem with the device.

Logic flow: The following rules cause this rule to be invoked:
 DAS100: Volume with worst overall performance

Discussion: CPEXpert checked PEND time, CONN time, DISC time, and IOSQ time. For legacy devices, CPEXpert also has checked seek time and missed RPS reconnect time. None of these areas **consistently** accounted for a majority of the delay for the volume selected as the volume with the "worst" overall performance. None of these areas accounted for the delay for the volume selected as the volume with the "worst" overall performance.

Suggestion: This condition most likely would occur when the guidance provided to CPEXpert is too restrictive. For example, if a large number of volumes are excluded from analysis (using the EXCLUDE option), the remaining volumes may have no particular problem. However, the logic is designed to select a "worst" volume, regardless of whether that volume actually has problems.

Rule DAS180: APPLICATIONS ACCESSING VOLUME WITH WORST PERFORMANCE

Finding: CPExpert identifies the applications accessing the volume with the worst performance.

Impact: This finding is used to select applications for potential rescheduling, file movement, or design changes. This rule is applicable only if CPExpert is performing "expanded" analysis of the data.

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with worst overall performance

Discussion: CPExpert lists the applications accessing the volume with the worst performance during the entire measurement period being analyzed. The list is ordered descendingly by I/O operations, so you can identify the applications that probably had the most impact on the volume. |

As described in Section 5, the SMF Type 30(Interval) information might not be synchronized with the SMF Type 70(series) information. In this case, the identification of applications may not correctly identify the applications with the most I/O operations to the volume. However, if the analysis produces the same results after analyzing more than one day's measurement data, you can be more comfortable that the applications are correctly identified. |

In some cases, the Type 30(Interval) data does not reflect all I/O activity to a device. This normally is caused by the Interval records not being written for all address spaces (such as started tasks running in SYSSTC service class). Consequently, you should use the output from DAS180 with care. |

The following example illustrates the output from Rule DAS180: |

RULE DAS180: APPLICATIONS ACCESSING THE VOLUME WITH WORST PERFORMANCE

The following applications, executing on system SYSF, accessed VOLSER SYF064 during the overall measurement interval (ranked by the number of I/O operations):

JOB	STEP	SERVICE CLASS	TYPE 30 START	INTERVAL END	EXCP RATE	AVG CONNECT TIME PER SEC
CXND9940	CXN99410	BATPMED.	0:00:00	1:00:00	203.3	0.269
CXND9940	CXN99410	BATPMED.	23:00:00	0:00:00	164.0	0.218
CXHD3015	CXH31550	BATPROD.	23:11:05	0:00:00	80.5	0.107
CXHD3015	CXH31550	BATPROD.	0:00:01	0:27:48	78.6	0.104
DB2HDBM1	GRP#PDBH	CICSHI..	7:00:00	8:00:00	0.3	0.000
DB2HDBM1	GRP#PDBH	CICSHI..	8:00:00	9:00:00	0.1	0.000
BC888222	BCVPGP15	BATPROD.	13:14:57	13:15:11	1.4	0.002
CXPD9042	CXP92000	BATPROD.	5:30:00	5:30:03	2.6	0.005
NALBM774	CXP12020	WLMLow..	3:29:50	3:30:08	0.3	0.001
N0932L0F	CXP12018	WLMLow..	4:45:00	4:45:05	1.2	0.002

Suggestion: In many cases, once the applications are identified, either applications personnel or systems personnel will verify whether the particular applications are responsible for the majority of the I/O operations to the particular volume.

Once you are comfortable that the applications are responsible for the majority of the I/O operations to the "worst" performing device, you can take action to (1) reschedule the applications to minimize contention, (2) examine the files referenced by the applications, or (3) change the way in which the applications access their files.

Rule DAS185: APPLICATIONS ACCESSING VOLUME DURING MEASUREMENT INTERVAL WITH WORST PERFORMANCE

Finding: CPEXpert identifies the applications accessing the volume during the the measurement interval with the worst performance.

Impact: This finding is used to select applications for potential rescheduling, file movement, or design changes. This rule is applicable only if CPEXpert is performing "expanded" analysis of the data.

Logic flow: The following rules cause this rule to be invoked:
 DAS100: Volume with worst overall performance

Discussion: CPEXpert lists the applications accessing the volume with the worst performance during the measurement period with the worst performance being analyzed. The list is ordered descendingly by I/O operations, so you can identify the applications that probably had the most impact on the volume.

As described in Section 5, the SMF Type 30(Interval) information might not be synchronized with the SMF Type 70(series) information. In this case, the identification of applications may not correctly identify the applications with the most I/O operations to the volume. However, if the analysis produces the same results after analyzing more than one day's measurement data, you can be more comfortable that the applications are correctly identified.

In some cases, the Type 30(Interval) data does not reflect all I/O activity to a device. This normally is caused by the Interval records not being written for all address spaces (such as started tasks running in SYSSTC service class). Consequently, you should use the output from DAS180 with care.

However, if the analysis produces the same results after analyzing more than one day's measurement data, you can be more comfortable that the applications are correctly identified.

Suggestion: In many cases, once the applications are identified, either applications personnel or systems personnel will verify whether the particular applications are responsible for the majority of the I/O operations to the particular volume.

Once you are comfortable that the applications are responsible for the majority of the I/O operations to the "worst" performing device, you can take

action to (1) reschedule the applications to minimize contention, (2) examine the files referenced by the applications, or (3) change the way in which the applications access their files.

Rule DAS200: VOLUME PROVIDING WORST OVERALL PERFORMANCE TO CRITICAL WORKLOAD

Finding: The identified volume provided the worst overall performance to the "loved one" workload during the entire measurement period. RULE DAS200 is quite similar to RULE DAS100; RULE DAS200 applies only to DASD devices accessed by critical (or "loved one") workload, while RULE DAS100 applies to **all** DASD devices.

Impact: This finding will have a HIGH IMPACT on the performance of the "loved one" workload.

Logic flow: This is a basic rule finding; there are no predecessor rules.

Discussion: CPExpert determines the average device response time, by device type, for each measurement interval. A "device type" for this purpose is any unique device type (e.g., IBM-3380 or IBM-3390), with the device type modified to reflect whether the device is cached, is a Parallel Access Volume (PAV), or is a paging device.

The purpose of determining the average device response time, by device type, is the underlying principle that there is little point in analyzing a particular device if its response time is better than average. Rather, the most improvement potential resides with devices whose response time is worse than average.

CPExpert selects a device in each measurement interval for further analysis if the device response time exceeds the average for its device type **and** the device was referenced by the "loved one" workload.

CPExpert consolidates information in various SMF records to build a model of the I/O configuration. This model includes utilization and queuing information for all channel paths, controllers, and devices. In creating the model, CPExpert:

- C Processes RMF Type 70 records to identify the systems that are in the sysplex.
- C Processes RMF Type 73 records to identify the physical channels that are associated with each system, and the type of channel (e.g., ESCON, FICON-Bridge, FICON-Native, etc.). Additionally, the physical and LPAR channel busy time is acquired for each system.

-
- C Processes RMF Type 78 records to obtain the logical control units associated with each system, and the channels associated with each logical control unit. Additionally, controller busy and director port busy times are acquired.
 - C Processes RMF Type 74 records to obtain devices associated with each logical control unit. Device performance characteristics are also acquired from the Type 74 records.

The result from the above is a record for each device, containing information about the devices; and the logical control units, channels, and systems that are associated with each device.

CPEXpert constructs a frequency distribution of all devices whose response is worse than average for the type of device, weighted by the number of I/O operations executed by the device. This yields a weighted measure of the potential performance improvement that might be achieved for each device. This frequency distribution is sorted descending, to yield an ordered list of the devices with the most improvement potential. This ordered list represents an "ordered intensity of access" distribution of the devices.

CPEXpert selects the top devices from the ordered list of devices referenced by the "loved one" workload. These represent the devices with the most improvement potential with respect to the "loved one" workload. CPEXpert reports information about the top devices from the list, by sysplex and by system (see Rule DAS000 and Rule DAS050 for additional information about this information).

Detailed information regarding the "worst" device is extracted and reported for each measurement interval. (RULE DAS290 reports summary results for the remaining top devices.)

Please refer to Section 5 for a discussion regarding some of the limitations of the technique of associating workload to device performance. In brief summary, the Type 30 information is not synchronized with the Type 70(series) information. Consequently, the analysis is approximate, rather than precise.

Exhibit DAS200-1 provides a sample output resulting from the analysis. The VOLSER and device number of the "worst" performing device are identified in the narrative. Information is provided about the overall average I/O rate and the device utilization for the entire measurement period being analyzed.

RULE DAS200: VOLUME WITH WORST OVERALL PERFORMANCE

VOLSER RSA002 (device 0194) provided CICS with the worst overall performance during the entire measurement period (13:59, 23JUL1991 to 16:59, 23JUL1991). This pack had an overall average of 17.2 I/O operations per second, was busy processing I/O for an average of 47% of the time, and had I/O operations queued for an average of 20% of the time. The following summarizes significant performance characteristics of VOLSER RSA002:

MEASUREMENT INTERVAL	I/O RATE	I/O RESP	--AVERAGE PER SECOND DELAYS--	MAJOR PROBLEM
			CONN DISC PEND IOSQ	
15:14-15:29,23JUL1991	16.9	0.665	0.096 0.288 0.073 0.208	QUEUING
15:29-15:44,23JUL1991	19.0	0.825	0.104 0.347 0.101 0.273	QUEUING
15:44-15:59,23JUL1991	18.5	0.737	0.103 0.350 0.053 0.231	QUEUING
15:59-16:14,23JUL1991	16.8	0.688	0.096 0.298 0.061 0.233	QUEUING
16:14-16:29,23JUL1991	19.9	0.793	0.108 0.354 0.058 0.273	QUEUING
16:29-16:44,23JUL1991	19.6	0.874	0.102 0.356 0.105 0.311	QUEUING
16:44-16:59,23JUL1991	17.5	0.669	0.098 0.320 0.061 0.190	QUEUING

VOLUME WITH WORST OVERALL PERFORMANCE

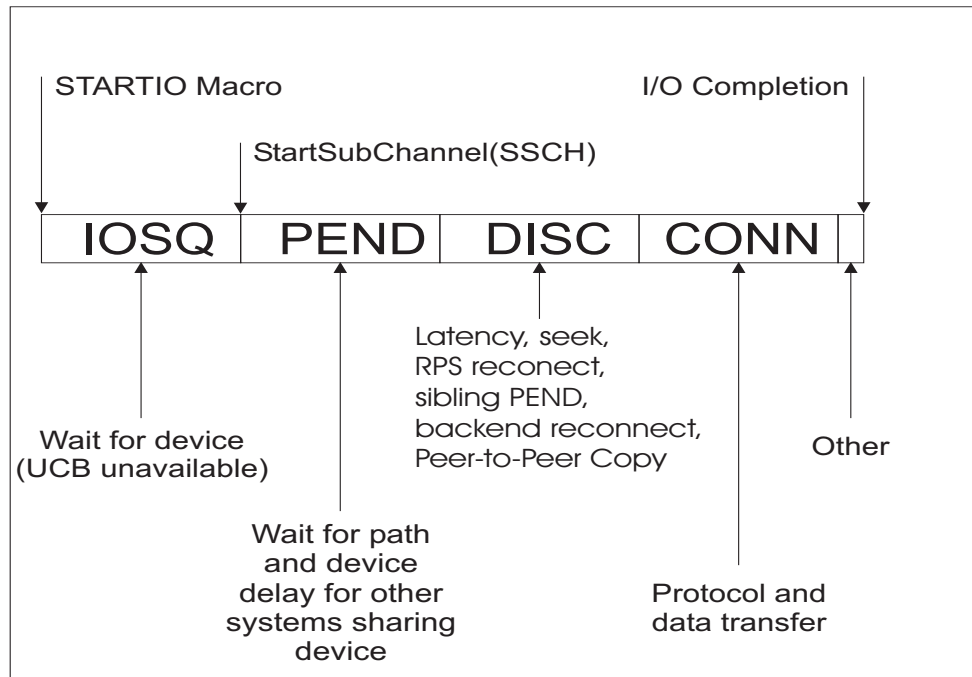
EXHIBIT DAS200-1

As shown in Exhibit DAS200-1, CPExpert provides a summary for each measurement interval, showing the average I/O rate for the interval, and the **average delay time per second** during the interval (the total is shown as **I/O RESP** in Exhibit DAS200-1). The average delay time per second effectively reflects the percent of each second (shown in milliseconds) in which the associated delay occurred.

CPExpert determines the major cause of device response delays (simply dividing each potential delay by the total device response time). The result from these calculations is evaluated to determine whether any area significantly predominates the I/O response time. If so, the respective area is listed as the major problem with the volume during the interval being analyzed.

From a high-level view, there are four key measures of DASD performance: IOS Queue (IOSQ) time, PENDING (PEND) time, disconnect (DISC) time, and connect (CONN) time. CPExpert provides information on these four key measures, and identifies the major cause of response delay.

The following figure illustrates these four measures and another potential element of DASD I/O time, titled "Other":



- C **IOSQ time.** IOSQ time is the time from the issuance of a STARTIO macro until the StartSubChannel (SSCH) instruction is issued. After the STARTIO macro is issued, the software determines whether the device is busy with *this system*; that is, whether there is an available Unit Control Block (UCB) for the device. If the device is not busy with *this system* (a UCB is available), the SSCH instruction is issued. However, if the device is busy with *this system*, the I/O request is queued. Thus, IOSQ time always means that the device is unable to handle additional requests from *this system*. (The emphasis on "this system" is explained in the below discussion of PEND time.)

This discussion of IOSQ time does not always apply to Parallel Access Volumes (PAVs)¹. With PAV devices, MVS creates multiple UCBs for each device, depending on how many "alias devices" have been defined. The multiple UCBs allow multiple active concurrent I/Os on a given device when the I/O requests originate from the same system². Using PAVs can dramatically improve I/O performance by nearly eliminating IOSQ.

¹PAV devices are available with Enterprise System Storage (ESS). With PAV devices, a "base device" address is defined, and a UCB is associated with this base address. "Alias device" addresses can be defined and UCBs are associated with the alias device addresses.

²Multiple Allegiance allows multiple active concurrent I/O operations on a given device when the I/O requests originate from different systems.

Please see Rule DAS250 for a more complete discussion of IOSQ time.

- C **PEND time.** PEND time is the time from the issuance of the StartSubChannel (SSCH) instruction until the device is selected by the control unit and physical positioning commands (such as seek and set sector) are transferred to the device. With modern fixed block architecture (FBA) devices, the PEND time ends when the physical positioning commands are presented to the *logical volume control block* within the control unit. The PEND time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, wait for device, or wait for “other” reasons)³.

The PEND time can be caused by the device being busy from *another system*. In this case, the system issuing the STARTIO macro (*this system*) would have no knowledge that the device was busy with another system. Rather, if a UCB were available for the device, the SSCH would be issued. However, the device could not necessarily be selected (unless multiple allegiance were available), since the device would be busy from another system.

Additionally, PEND time could accumulate even with PAV devices if the access were to an extent that was busy with another I/O operation from *this system*.

Please see Rule DAS230 for a more complete discussion of PEND time.

- C **DISC time.** DISC means that there is some delay that is often (but not always) associated with a mechanical movement during which the device disconnects from the control unit.

With legacy systems (e.g., 3380 drives attached to 3990-2 control units), the DISC time of most concern was associated with seek (arm movement) and rotational position sensing (time waiting for the disk platter to rotate to the location where desired data resides). Considerable performance improvement efforts were directed at reducing the seek activity and reducing the rotational position sensing

³PEND time is significantly reduced with FICON channels. FICON channels can have multiple I/O operations concurrently active, which reduces the potential PEND time caused by channel busy. There is no port busy time with FICON switches, and control unit time is significantly reduced. This statement regarding PEND time is not necessarily correct if a large number (more than 5) I/O operations are concurrently executing on a FICON channel. Dr. H. Pat Artis and Mr. Robert Ross have presented the results of research indicating that performance degrades significantly when more than 5 I/O operations (Open Exchanges) are concurrently active on a FICON channel (see “Understanding FICON Channel Path Metrics” at www.perfassoc.com).

(RPS)⁴ delays for the legacy systems. These two mechanical delays still exist for most modern *redundant array of independent disks* (RAID)⁵ systems, but their impact can not be directly reduced with normal methods.

With modern disks, data is cached into Actuator Level Buffers (ALBs), that contain data read from a track on the disk platter. Using ALBs can eliminate the RPS delays for records read on a particular track, since required data is read into the device buffer during a single rotation and stored until a path is available to transfer the data. However, if a record is to be read from a new track, some RPS delay could exist since the record would not be in the ALB, and must be read from the new track. Some initial RPS delay would apply in this case. This initial RPS delay is neither measured nor preventable.

Additionally, data is cached into increasingly large cache on the controller. For a read operation, desired data often is found in the cache. Write operations normally end as the data to be written is placed in non-volatile storage (NVS); and the storage processor writes the data to the device asynchronous with other activity (as a “back end” staging operation).

Consequently, DISC time for modern systems is a result of *cache read miss* operations, potentially back-end staging delay for write operations, peer-to-peer remote copy (PPRC) operations, and other miscellaneous reasons⁶. DISC time often can be very small with adequate cache. For example, there would be zero disconnect time for a cache read hit (the record was found in the cache).

Please see Rule DAS260 for a more complete discussion of DISC time.

- C **CONN time.** CONN time includes the data transfer time, but also includes protocol exchange⁷ (or “hand shaking”) between the various components at several stages of the I/O operation.

⁴RPS delays are caused by a path not being available when the required data came under a device read head. Since a path was not available, the data could not be read and another rotation of the platter was experienced until the data again came under the device read head. Multiple rotations might be required, depending on the busy level of the path.

⁵An array is an ordered collection of physical devices (disk drive modules) that are used to define logical volumes or devices.

⁶Artis has described a “sibling PEND” condition that results from collisions within the physical disk subsystem of RAID devices. See “Sibling PEND: Like a Wheel within a Wheel,” www.cmg.org/cmgpap/int449.pdf.

⁷Note that the protocol exchange occurs at multiple points in the normal I/O operation, even though it is shown only once in this exhibit.

For devices attached to paths that include parallel channels and ECON channels, the data transfer time is simply the number of bytes transferred divided by the transfer speed. This is because a parallel channel or ESCON channel can have only one data transfer operation in execution at one time.

For devices attached to paths that include FICON channels, the algorithm is more complicated. This primarily is because a FICON channel can perform multiple data transfer (read and write) operations at one time. The data packets for multiple read or write operations are interleaved (or multiplexed) in the FICON link. CONN time for an individual I/O begins with the first frame of data transferred and ends last frame of data transfer, even though data for other I/O operations might be transferred concurrently on the link. Consequently, if multiple data packets (representing data for multiple read or write operations) are interleaved on the FICON link, the elapsed time for any particular I/O operation can be elongated⁸ when compared with the elapsed time of the same I/O operation on an ESCON channel.

Please see Rule DAS240 for a more complete discussion of CONN time.

C **OTHER time.** There are at least two other potential I/O delays for DASD: (1) waiting for the I/O completion interrupt to be serviced by a processor and (2) waiting for the I/O interrupt to be serviced by a domain under PR/SM. Neither potential I/O delay is expected to be of the magnitude of the four "standard" I/O delays. However, they can be significant in special circumstances.

C Multi-processor configurations can use any processor to service an I/O interrupt. However, when a processor services an I/O interrupt, the processor's high-speed cache storage is no longer valid when control is returned to the interrupted task. Consequently, many of the processor's high-performance design features may be nullified.

A hardware feature allows processors to be disabled for I/O interrupts. With this method, only a small number (perhaps only one) processor is enabled for interrupt processing. Only this processor will have its high-speed cache storage disturbed by the task-switching required for interrupt processing, and only this

⁸The relative speed of a FICON channel is much higher than that of an ESCON channel. Consequently, the elapsed time of any particular I/O operation should be less on a FICON channel than on an ESCON channel, even if there are multiple I/O operations interleaving data. This statement regarding elapsed time is not necessarily correct if a large number (more than 5) I/O operations are concurrently executing on a FICON channel. Dr. H. Pat Artis and Mr. Robert Ross have presented the results of research indicating that performance degrades significantly when more than 5 I/O operations (Open Exchanges) are concurrently active on a FICON channel (see "Understanding FICON Channel Path Metrics" at www.perfassoc.com).

processor will periodically have its high-performance design features nullified. The disadvantage to this approach is that an interrupt may occur while the processor is busy servicing a previous interrupt.

If an interrupt is pending and no processor is enabled to service the interrupt, the interrupt must wait until a processor is available. This time should be insignificant, unless the system is processing a significantly large number of I/O operations. If the system is processing a large number of I/O operations (or if the I/O is particularly time-sensitive), the interrupt pending delay could pose performance problems.

After the processor completes processing for an I/O interrupt, it issues a Test Pending Interrupt (TPI) instruction to determine whether there are any interrupts pending. If an I/O interrupt is pending, the processor proceeds to service that interrupt.

The CPENABLE keyword in the IEAOPTxx member of SYS1.PARMLIB is used to specify the percent of I/O interrupts detected by the TPI instruction, compared with all I/O interrupts. When the percent exceeds the high threshold of the CPENABLE keyword, MVS enables another processor to handle pending I/O interrupts. If the percent falls below the low threshold of the CPENABLE keyword, MVS will disable a processor (to the point that only one processor is enabled). IBM's recommended setting for the CPENABLE keyword differs, depending on the level of processor.

- C MVS environments running under as a guest under VM or in a logical partition (LPAR) under PR/SM are subject to I/O interrupt delays. These delays can occur if another guest (for VM) or another domain is in its dispatch interval when the I/O interrupt completion is posted. The I/O interrupt remains pending until the guest or domain is dispatched. These delays have been estimated to be far more significant than might otherwise be expected.

Suggestion: There are no suggestions directly associated with this rule. Subsequent rules will analyze the device problems and attempt to determine the cause of poor performance.

Rule DAS202: VOLUME PROVIDING NEXT WORST OVERALL PERFORMANCE TO CRITICAL WORKLOAD

Finding: The identified volume provided the next worst overall performance to the "loved one" workload during the entire measurement period.

Impact: This finding will have a HIGH IMPACT on the performance of the "loved one" workload.

Logic flow: This is a basic rule finding; there are no predecessor rules.

Discussion: Rule DAS202 identified the volume that had the worst overall performance during the entire measurement period, from the perspective of the "loved one" workload. Rule DAS202 is produced for each successive "worst performing" device selected from an ordered list (Rule DAS290 shows the ordered list of devices).

The number of devices analyzed by Rule DAS202 (and successive rules resulting from the analysis of each device) is controlled by the **ANALYZE** guidance variable (see Section 3: Specifying Guidance Variables).

Please refer to Rule DAS202 for a discussion of the information presented with Rule DAS202.

Suggestion: There are no suggestions directly associated with this rule. Subsequent rules will analyze the device problems and attempt to determine the cause of poor performance.

Rule DAS205: VOLUME PERFORMANCE WAS NOT CONSISTENTLY POOR

Finding: The performance of the volume having the "worst" overall performance for the critical workload was not consistently poor in any single area.

Impact: This finding has little impact, other than to indicate that there may be no problem with the device.

Logic flow: The following rules cause this rule to be invoked:
DAS200: Volume with worst overall performance

Discussion: CPEXpert checked PEND time, CONN time, DISC time, and IOSQ time. For legacy devices, CPEXpert also has checked seek time and missed RPS reconnect time. None of these areas **consistently** accounted for a majority of the delay for the volume selected as the volume with the "worst" overall performance with respect to the critical workload.

Suggestion: This condition most likely would occur when the period(s) of poor performance were extremely poor or if the guidance provided to CPEXpert is too restrictive.

- Devices sometimes have periods of extremely poor performance combined with a large number of I/O operations. This combination may result in the device being selected as the one with the most potential for performance improvement, even though the device did not have consistently poor performance. The most likely cause of such a situation is a particular application or combination of applications executing in the interval with poor performance.
- The guidance provided to CPEXpert may be too restrictive. For example, if a large number of volumes are excluded from analysis (using the EXCLUDE option), the remaining volumes may have no particular problem. However, the logic is designed to select a "worst" volume, regardless of whether that volume actually has problems.

Rule DAS210: SEEKING WAS THE MAJOR CAUSE OF RESPONSE DELAY

Finding: Seeking was the major cause of the I/O response delay with the device.

Impact: This finding can have a MEDIUM IMPACT or HIGH IMPACT. Since the device is the "worst performing" device for the critical (or "loved one" workload, this finding can have a HIGH IMPACT, depending upon the amount of seeking being done. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).*

Logic flow: The following rule causes this rule to be invoked:
 DAS200: Volume with the worst overall performance

Discussion: The discussion associated with Rule DAS200 describes how CPExpert creates a model of the I/O configuration based upon the IOCP macros and RMF data.

CPExpert applies queuing formulae to the model to estimate the amount of delay attributed to missed RPS reconnect (these delays are a function of the probability of a device finding all paths busy when the device tries to reconnect to the channel path).

The estimated missed RPS reconnect time is subtracted from the DISC time reported by RMF. Additionally, the average latency for the device type is subtracted from the DISC time. The resulting time is assumed to be the seek time. (Note the below discussion about why this assumption might not be correct.)

CPExpert performs the above analysis for each measurement interval reflected in the data. Rule DAS210 is produced if seeking was the major problem for a majority of the measurement intervals.

There are potential problems with this approach, although the approach is generally used throughout the computer industry as a way of estimating missed RPS delays and of estimating seeking.

- The queuing formulae assume exponential interarrival times, exponential service distributions, and an infinite population (the M/M/c formula - Erlang's C formula - is used for the calculations). These assumptions may not be correct if, for example, the I/O activity is a function of a single application.

In his class "MVS I/O Configuration Management", Dr. Jeffery Buzen provides a Dump/Restore application as an excellent example of an application that does not follow standard queuing assumptions.

- The device may be cached, and it may be impossible to apportion the DISC time residual after subtracting missed RPS reconnect time. This time may represent a few missed cache read operations with long seek distances, or may represent a relatively large number missed cache read operations with little seeking but the standard latency for the device.

Thus, the seeking analysis can only show potential problems, rather be considered a definitive indication. However, it is usually a fairly accurate indication of the problem. If high average seeking is reported, you can be fairly certain that high seeking did occur. This is particularly true if the problem is reported throughout the measurement intervals. The "uncertainty" tends to be related to relatively low seeking or seeking reported for cached devices.

Rule DAS210 reports the overall average number of milliseconds out of each second in which the device was positioning the arm **while servicing I/O requests for the "loved one" workload**. It is important to remember that the device may have had different performance characteristics at different measurement intervals when the device was not used by the "loved one" workload.

Additionally, Rule DAS210 summarizes key information about the period of worst performance, if seeking was the major cause of delay during this period.

Suggestion: The seeks can be minimized by (1) rearranging files within the pack, (2) moving files from the pack to another actuator, (3) changing the application file accessing characteristics, or (4) possibly restricting the applications allowed to access the pack.

Rule DAS220: MISSED RPS RECONNECT WAS MAJOR CAUSE OF I/O DELAY TO CRITICAL WORKLOAD

Finding: CPExpert has determined that missed Rotational Position Sensing (RPS) reconnects was a major cause of delay in DASD response for the device.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).*

Logic flow: The following rules cause this rule to be invoked:
DAS200: Volume with the worst overall performance

Discussion: CPExpert computes the average channel path busy for paths to the device. SMF Type 78 information is used to acquire the channel path utilization, and the IOCP macro information is used to determine the channel paths to the device.

If the device is capable of DLS mode, two channel paths can be concurrently busy to the string. If the device is capable of DLSE mode, four channel paths can be concurrently busy to the string.

CPExpert use an M/M/C queuing model (Erlang's C formula) to compute the probability that all paths to the device were busy when the device attempted to reconnect. The M/M/C queuing model is adjusted depending upon the number of concurrent paths. (Please refer to *Probability, Statistics, and Queuing Theory* by Arnold O. Allen for a description of the M/M/C queuing model.)

The average number of missed RPS reconnect attempts is given by the formula:

$$N = \frac{U}{1 - U}$$

where U is the probability that a reconnect attempt will find all paths busy

The average time spent attempting to reconnect is simply the rotational time of the device multiplied by the average number of missed RPS reconnect

attempts. This yields an estimate of the average I/O delay caused by missed RPS reconnect attempts.

Rule DAS220 is produced if the average time spent attempting to reconnect to the channel path accounts for a significant percent of the device response time **to the "loved one" workload**.

The I/O delay caused by missed RPS reconnect may be under-estimated if the device is shared between systems. This is because the controller may be busy to a different system when the device attempts to reconnect. There is no information in RMF records to show the time the controller was busy to another system. RMF data from all systems sharing the controller must be consolidated and analyzed in order to obtain this information. This analysis is performed by the DASD Component only if you have selected the "analysis of shared DASD conflicts" option (see Section 3, Chapter 2.5).

Suggestion: Missed RPS reconnect time is caused by too much activity on the channel paths to the device. This problem can be corrected by:

- Removing I/O data transfer from the path(s). When devices are seeking or searching for a sector, the channel path is not busy. Therefore, high channel activity is primarily due to transferring data (controller protocol accounts for some path busy, but this time generally is quite small). Removing I/O data transfer from the path(s) can be accomplished by:
 - Moving data sets from the paths. This often is the easiest solution, since data sets with high data transfers can be relocated to other volumes on different channel paths.
 - The volumes responsible for high data transfer time could be moved to other strings on different channel paths.
- Rescheduling workloads to minimize the contention. For example, some batch jobs may be performing heavy I/O to volumes on the string. These jobs may be rescheduled to a time when their I/O would not cause problems.
- Adding paths to the device. If the device is not capable of DLS or DLSE mode, this will require upgrading the device (e.g., from IBM-3380 Model A04 to a more recent model). If the device is capable of DLS mode (that is, it can dynamically reconnect to either of two paths), you may consider upgrading the device to one capable of DLSE mode (that is, it can dynamically reconnect to any of four paths). This action may require that you upgrade your controller from an IBM-3880 storage controller to an IBM-3990 storage controller. Since upgrades are relatively expensive,

you should first assess the feasibility of moving active data sets or volumes from the path.

- Adding channel paths, acquiring additional controllers, and moving some volumes to the new controllers. This would be a more expensive option, and may not be feasible (depending upon the processor and I/O configuration).

Rule DAS221: VOLUMES CONTRIBUTING TO RPS DELAY

Finding: CPExpert identifies the volumes contributing to missed RPS reconnect delays to the "loved one" workload.

Impact: This information can be useful when deciding on a course of action to correct the missed RPS reconnect problems. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).* |

Logic flow: The following rules cause this rule to be invoked:
DAS200: Volume with the worst overall performance
DAS220: Missed RPS reconnect was major cause of I/O delay

Discussion: If Rule DAS220 is produced, CPExpert examines the SMF Type 74 information and IOCP macro information to select all devices sharing paths with the device experiencing missed RPS reconnect delays. Up to 10 of these devices are listed, ordered descendingly by their contribution to path utilization.

A device contributes to path utilization mostly based upon the connect time of the device. Consequently, CPExpert displays the devices having the largest per-second connect times on the path.

Suggestion: You should consider separating the volumes contributing to missed RPS delays from those on the volume experiencing the significant missed RPS delays.

- You can move data sets from the volumes contributing the most to path utilization.
- You can move volumes to a different string.
- You might reschedule workload to minimize the path utilization at critical times.

Rule DAS223: NON-DASD DEVICES CONTRIBUTED TO RPS DELAY

Finding: CPExpert has determined that non-DASD I/O devices were attached to a channel path of the volume experiencing missed RPS reconnect delays. These non-DASD I/O devices were busy a significant percent of the time and contributed to the RPS delay.

Impact: This finding can have a HIGH IMPACT on the performance of the device experience the missed RPS reconnect delays. Since this device is used by the "loved one" workload and is the volume with the worst overall performance, this finding can have a HIGH IMPACT on the performance of the "loved one" workload. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).*

Logic flow: The following rules cause this rule to be invoked
DAS200: Volume with the worst overall performance
DAS220: Missed RPS reconnect was major cause of I/O delay

Discussion: CPExpert determines whether any non-DASD I/O devices (e.g., tapes drives, etc.) share channel paths with DASD. If missed RPS reconnect delays were a major cause of I/O delay, CPExpert undertakes an analysis of the non-DASD I/O devices sharing channel paths. CPExpert examines the SMF Type 74 information to determine whether these non-DASD devices had a significant connect time to the path.

CPExpert uses a M/M/c queuing model to estimate the amount of missed RPS reconnect delay caused by the path utilization of the non-DASD devices.

Rule DAS223 is produced if the queuing model estimates that path utilization of the non-DASD devices causes more than 10% of the missed RPS reconnect delay.

Suggestion: CPExpert suggests that you eliminate or minimize the impact of the non-DASD I/O devices on the DASD performance by considering the following alternatives:

- Reschedule the workload accessing the non-DASD I/O devices to a period when the data transfer would not cause DASD problems.

-
- Remove the non-DASD I/O devices from the channel paths serving the DASD devices. This may mean that you must acquire additional channel paths.
 - If neither of the above options are feasible, consider placing only low-utilization (and non-critical) DASD on the paths shared with the non-DASD I/O devices.

Rule DAS225: APPLICATIONS CONTRIBUTING TO RPS DELAY

Finding: CPExpert identifies the applications (other than the "loved one" workload) contributing to missed RPS reconnect delays.

Impact: This information can be useful when deciding on a course of action to correct the missed RPS reconnect problems. *This finding applies only to legacy systems (e.g., 3380 devices attached to 3990-2 controllers).* |

Logic flow: The following rules cause this rule to be invoked:
DAS200: Volume with the worst overall performance
DAS220: Missed RPS reconnect was major cause of I/O delay

Discussion: If Rule DAS220 is produced, CPExpert examines the SMF Type 30(Interval) information to select all applications (other than the "loved one" workload) that reference the volume. An application contributes to path utilization (and thus, contributes to missed RPS reconnect) mostly based upon the connect time of I/O operations to the device.

CPExpert lists the applications accessing the volume with the worst performance for the "loved one" workload, during the entire measurement period being analyzed. The list is ordered descendingly by the average percent use of paths, so you can identify the applications that probably had the most impact on the volume.

As described in Section 5, the SMF Type 30(Interval) information is not synchronized with the SMF Type 70(series) information. Consequently, the identification of applications may not correctly identify the applications with the most I/O operations to the volume.

However, if the analysis produces the same results after analyzing more than one day's measurement data, you can be more comfortable that the applications are correctly identified.

Suggestion: In many cases, once the applications are identified, either applications personnel or systems personnel will verify whether the particular applications are responsible for the majority of the I/O operations to the particular volume.

You should become comfortable that the applications presented are responsible for the utilization of the paths during the period when missed RPS reconnect delayed the "loved one" workload. You can then take action

to (1) reschedule the applications to minimize contention, (2) examine the files referenced by the applications, or (3) change the way in which the applications access their files.

Rule DAS230: PEND TIME WAS MAJOR CAUSE OF I/O DELAY

Finding: CPExpert has determined that excessive PEND time was a major cause of delay in DASD response for the device.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
 DAS200: Volume with the worst overall performance

Discussion: PEND time is the time from the issuance of the StartSubChannel (SSCH) instruction until the device is selected by the control unit and physical positioning commands (such as seek and set sector, or define extent) are transferred to the device.

With modern fixed block architecture (FBA) devices, the PEND time ends when the physical positioning commands are presented to the *logical volume control block* within the control unit. The PEND time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, or wait for device, or wait for “other” reasons)¹.

PEND is measured by the channel subsystem. After IOS issues the Start Subchannel command, the channel subsystem may not be able to initiate the I/O operation if any path or device busy condition is encountered:

C The channel selected for the I/O operation could be busy with another I/O operation from another system image in the same CEC. This time is not reflected in the SMF data.

C The director port could be busy with another I/O operation². This time is reflected in SMF data as SMF74DPB.

¹PEND time is significantly reduced with FICON channels. FICON channels can have multiple I/O operations concurrently active, which reduces the potential PEND time caused by channel busy. There is no port busy time with FICON switches, and control unit time is significantly reduced. This statement regarding PEND time is not necessarily correct if a large number (more than 5) I/O operations are concurrently executing on a FICON channel. Dr. H. Pat Artis and Mr. Robert Ross have presented the results of research indicating that performance can degrade significantly when more than 5 I/O operations (Open Exchanges) are concurrently active on a FICON channel (see “Understanding FICON Channel Path Metrics” at www.perfassoc.com).

²Director port busy can occur only on an ESCON channel. The use of buffer credits on a FICON native channel eliminates director port busy.

C The control unit could be busy with another I/O operation from another system. This time is reflected in SMF data as SMF74CUB.

C The device could busy with I/O from another system. This time is reflected in SMF data as SMF74DVB.

There can be “other” PEND time not reflected in the above descriptions. For many systems, “other” PEND time is zero or very small. For some systems, the “other” PEND time is dramatically large (often, 75% or more of the average response time).

Suggestion: Please refer to Rule DAS130 for further information about PEND time.

Rule DAS231: PEND DELAY TIME WAS CAUSED BY CHANNEL ACTIVITY

Finding: CPEXpert has determined that the excessive PEND time was caused by utilization of the channels to the device.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
DAS200: Volume with the worst overall performance
DAS230: Major cause of I/O delay was PEND time

Discussion: PEND time is the time from the issuance of the SSCH instruction until the device is selected by the control unit. This time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, wait for device, or wait for other reasons).

With modern fixed block architecture (FBA) devices, the PEND time ends when the physical positioning commands are presented to the *logical volume control block* within the control unit. The PEND time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, or wait for device, or wait for “other” reasons)¹.

The PEND time due to channel busy means the channel that was selected for the I/O operation was busy with another I/O operation from another system image in the same CEC. Since the channel was busy, the SSCH instruction could not result in device selection, and a PEND condition existed.

SMF Type 74 records do not contain the PEND time caused by channel busy². However, CPEXpert calculates an estimated PEND for channel busy based on I/O configuration information.

¹PEND time is significantly reduced with FICON channels. FICON channels can have multiple I/O operations concurrently active, which reduces the potential PEND time caused by channel busy. There is no port busy time with FICON switches, and control unit time is significantly reduced. This statement regarding PEND time is not necessarily correct if a large number (more than 5) I/O operations are concurrently executing on a FICON channel. Dr. H. Pat Artis and Mr. Robert Ross have presented the results of research indicating that performance can degrade significantly when more than 5 I/O operations (Open Exchanges) are concurrently active on a FICON channel (see “Understanding FICON Channel Path Metrics” at www.perfassoc.com).

²MXG contains a variable AVGPNDCHA, which is titled 'AVG (MS)*PEND DUE TO*CHANNEL BUSY'. However, the AVGPNDCHA variable is simply created from the AVGPNDIR titled 'AVG (MS)*PEND DUE TO*DIRECTOR PORT' variable. MICS does not contain a “PEND CAUSED BY CHANNEL BUSY” variable.

When CPExpert creates the model of the I/O configuration, it retains information about each path to a device. Included in this path information is the physical path busy at the CEC level, for each path. Consequently, CPExpert has an overall view of all physical paths to the device, and can calculate overall channel activity for all channels to the device.

Rule DAS231 is produced when the calculated PEND time due to channel busy accounts for more than one-third of the device PEND time. This output will be produced for all channels sets (by CEC serial number) that are used to reference the logical volume experiencing high PEND time.

The following example illustrates the output from Rule DAS231:

```

RULE DAS231:  LARGE PEND TIME DELAY WAS CAUSED BY CHANNEL BUSY.

A significant amount of the PEND time delay was caused by high channel
utilization for the channels connected to VOLSER CICS11.  This volume
was referenced by the indicated channels.

                                AVERAGE PHYSICAL CHANNEL BUSY FOR CHPIDS:
MEASUREMENT INTERVAL          3D  59  67  72  7F  99  A7  B4
8:30- 8:45,22OCT2001          19  18  84  33  53  42  31  23

```

Suggestion: If an important device is experiencing delays because of high PEND caused by channel utilization, you should consider the following alternatives:

C **Redistribute data sets.** The high PEND time might be solved by redistributing high activity data sets among different volumes on different paths.

If SMF Type 42 (Data Set Statistics) are available, CPExpert will identify data sets on the logical volume that have heavy I/O activity. However, please keep in mind that the PEND time is caused by channel activity. The I/O activity of the particular volume experiencing high PEND time might not be (and probably is not) the cause of high PEND time. Consequently, examining the results for the Type 42 (Data Set Statistics) for the volume might not yield satisfactory results.

C Move the logical volume to a different controller referenced by different channels.

C If redistributing the data sets or moving the volume is not feasible, perhaps more channels can be assigned to the logical control unit through which the device is referenced.

Rule DAS232: PEND DELAY TIME WAS CAUSED BY DIRECTOR PORT

Finding: A significant amount of the PEND time delay was caused by director port busy.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
 DAS200: Volume with the worst overall performance
 DAS230: Major cause of I/O delay was PEND time

Discussion: PEND time is the time from the issuance of the SSCH instruction until the device is selected by the control unit. This time is caused by queuing for the path (wait for channel, wait for control unit or wait for head-of-string), and can be caused by other systems sharing the device (wait for device).

The PEND time due to director port busy, means the director port was busy with another I/O operation from another system image in the same CEC. Since the director port was busy, the SSCH instruction could not result in device selection, and a PEND condition existed.

SMF Type 74 records contain the PEND time caused by director port busy (variable SMF74DPB). This variable is contained in MXG as AVGPNDIR and in MICS as DBADPBTM.

CPEXpert produces Rule DAS130 to report the causes of PEND delay time. Rule DAS232 is produced when director port busy delay was the major cause of PEND delay time.

Suggestion: If Rule DAS232 is consistently produced, you should consider the following alternatives:

C Redistribute data sets. The high PEND time might be solved by redistributing high activity data sets among different volumes on different paths.

If SMF Type 42 (Data Set Statistics) are available, CPEXpert will identify data sets on the logical volume that have heavy I/O activity. However, please keep in mind that the PEND time is caused by director port activity. The I/O activity of the particular volume experiencing high PEND

time might not be (and probably is not) the cause of high PEND time. Consequently, examining the results for the Type 42 (Data Set Statistics) for the volume might not yield satisfactory results.

C Move the logical volume to a different controller referenced by different channels.

C If redistributing the data sets or moving the volume is not feasible, then perhaps more channels can be assigned to the logical control unit through which the device is referenced.

Rule DAS233: PEND DELAY TIME WAS CAUSED BY CONTROLLER BUSY

Finding: A significant amount of the PEND time delay was caused by controller activity.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
 DAS200: Volume with the worst overall performance
 DAS230: Major cause of I/O delay was PEND time

Discussion: PEND time is the time from the issuance of the SSCH instruction until the device is selected by the control unit. This time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, wait for device, or wait for other reasons).

Large PEND times for devices attached to cached controllers may imply a high percent of read miss operations, or non-volatile storage (NVS) writes.

To improve the probability of a read hit, the controller can *prestage* data into its cache. Prestaging means that data is read into the controller's cache ahead of its actually being required for use by an application. The amount of data that is prestaged depends on (1) whether the data is being accessed in a direct (random) mode or in a sequential mode and (2) the controller model and the enhancements made to the controller.

C For *direct mode*, after the record is located, the 3390-3 and 3990-6 (initial version) stages in the balance of the track being read.

The 3990 Model 6 (with record cache) stages only the records requested into cache, eliminating the balance of the track staging that is normal with track caching as was implemented on initial versions of 3990-6 and on the 3990-3. This improvement reduces the PEND time caused by the controller busy during track staging.

C As examples of prestaging for *sequential mode*, the 3990-3 reads up to two tracks into the cache¹ before they are required, while the ESS 2105 sequential staging reads up to two cylinders ahead.

¹With the Sequential Staging Performance Enhancement, the 3990-3 can prestage up to a full cylinder (15 tracks) into the cache.

During prestaging operations for sequential reads, the control unit regularly checks to see whether other I/O requests are waiting to be processed. If any are waiting, the control unit interrupts the prestage operation, processes the queued requests, and continues with the prestage.

In DASD Fast Write Mode, the data is stored simultaneously in cache storage and in nonvolatile storage (NVS). At some subsequent time, the data in NVS can be *destaged* to DASD.

In Cache Fast Write Mode, data is placed into cache immediately, and there is no interaction with the device nor with NVS. However, if cache memory is required (or if Cache Fast Write Mode is turned off), the data in cache is destaged to DASD.

Significant PEND time can result from destaging to DASD.

Rule DAS234: PEND DELAY TIME WAS CAUSED BY DEVICE BUSY

Finding: A significant amount of the PEND time delay was caused by device busy.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
 DAS200: Volume with the worst overall performance
 DAS230: Major cause of I/O delay was PEND time

Discussion: PEND time is the time from the issuance of the SSCH instruction until the device is selected by the control unit. This time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, wait for device, or wait for other reasons). Some of the causes of PEND for device busy are listed below:

C PEND time for device busy can be caused by other systems in the sysplex that issue a RESERVE for the device. While the RESERVE is held, I/O operations to the device will be held in a PEND for device busy state.

C Multiple Allegiance allows multiple active concurrent I/O operations on a particular device when the I/O requests originate from different systems. With Multiple Allegiance, there is complete access with read I/O operations. For write I/O operations, there is concurrent access unless there is a conflicting extent¹. If there is a conflicting extent, the controller holds the I/O operation in a PEND state for the device.

C After an I/O operation, the device will read the remainder of the track into its device-level buffer. This is done to prevent delay for rotational positioning. If a new I/O operation is attempted while data is being read into the device cache buffer, the I/O operation will be in a PEND state for device busy.

Suggestion: If Rule DAS234 is produced frequently, you should consider the following alternatives.

¹ A conflicting extent is one in which the write operation attempts to update an extent.

C If shared device analysis is specified as **%LET SHARED = Y**; in **USOURCE(DASGUIDE)**, CPExpert will analyze potential problems caused by sharing DASD. Rule DAS300 will be produced for all systems that share the device with high PEND time, if CPExpert concludes that other systems could cause performance problems. The RESERVE time will be included in the output from Rule DAS300.

If this RESERVE time is high for the device, you should consider whether high activity data sets can be moved among different volumes on different paths.

C Alternatively, determine whether the data sets can be moved to a controller that supports Multiple Allegiance². Multiple Allegiance allows multiple active concurrent I/O operations on a particular device when the I/O requests originate from different systems. With Multiple Allegiance, there is complete access with read I/O operations. For write I/O operations, there is concurrent access unless there is a conflicting extent.

C Alternatively, consider whether workload scheduling can eliminate the conflicts between the data access requirements between systems.

²Multiple Allegiance is available with IBM's Enterprise Storage Server (ESS) subsystems.

Rule DAS235: PEND DELAY TIME WAS CAUSED BY OTHER DELAYS

Finding: A significant amount of the PEND time delay was caused by delays that were not unexplained by either the causes of PEND time as reported by SMF, or estimated causes as calculated by CPEXpert.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
 DAS200: Volume with the worst overall performance
 DAS230: Major cause of I/O delay was PEND time

Discussion: PEND time is the time from the issuance of the SSCH instruction until the device is selected by the control unit. This time is caused by queuing for the path (wait for channel, wait for director port, wait for control unit, wait for device, or wait for other reasons).

SMF contains information describing the PEND delay time caused by wait for director port (SMF74DPB), wait for control unit (SMF74CUB), and wait for device (SMF74DVB). As described in Rule DAS131, CPEXpert computes an estimated PEND time caused by channel busy.

After computing an estimated PEND for channel busy, CPEXpert subtracts the wait for channel, wait for director port, wait for control unit, and wait for device from the total Device PEND time reported by SMF in SMF74PEN. After this subtraction, there often is a remainder. This remainder¹ has been titled "other" PEND time since it is PEND time that is not reflected in the measured or estimated causes of PEND delay time.

For many systems, "other" PEND time is zero or very small. For some systems, the "other" PEND time is dramatically large (often, 75% or more of the average response time).

Suggestion: There are no suggestions with Rule DAS235. The finding is produced simply to alert you that the PEND delay data contains a large amount of

¹MXG does not calculate the PEND due to channel busy, but simply sets the AVGPNCCHA equal to the AVGPNDIR. MXG does subtract the PEND for director port, PEND for control unit, and PEND for device from the total device PEND delay. Since MXG does not calculate an estimated PEND time due to channel busy, MXG produces a much larger PEND "other" value than does CPEXpert's calculations.

delay that was unexplained by either the causes of PEND time as reported by SMF, or estimated causes as calculated by CPExpert.

At present, there is only conjecture² about additional cause of this “other” PEND time. Perhaps either IBM will better describe this “other” PEND time in future, or perhaps research will reveal likely causes of the “other” PEND time.

²According to MXG (ADOC74 comments), Dr. H. Pat Artis believes that the “other” PEND is often the internal response time of the subsystem, i.e., the time it takes the subsystem to accept, validate, and acknowledge the first Channel Control Word (CCS) of the channel program.

Rule DAS240: CONNECT TIME WAS A MAJOR CAUSE OF I/O DELAY

Finding: Connect time was a major cause of the I/O delay with the volume.

Impact: This finding may have a LOW IMPACT or MEDIUM IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
DAS200: Volume with the worst overall performance

Discussion: Connect time is the time in which the device is actually connected to the path. This time includes the data transfer time, but also includes protocol exchange (or "hand shaking") between the various components at several stages of the I/O operation.

The data transfer time obviously is a function of the amount of data being transferred. This simply is the number of bytes transferred divided by the transfer speed (for example, if 4096 bytes were transferred from an IBM-3380 with a transfer speed of 3,000,000 bytes per second, the 4096 bytes would require $4096/3,000,000$ seconds; or about 1.36 milliseconds).

Large connect times generally are caused by the following situations:

- A large average block size. This situation may be highly desirable for sequential data sets, but would be undesirable for randomly accessed data.
- Long multi-track searches. For example, the catalog must be searched for cataloged files, the Volume Table of Contents (VTOC) must be searched to find a requested file, a directory must be searched for partitioned data sets, etc.. These searches will result in long connect times for the volume involved.
- Program loading from system packs.

Rule DAS250: QUEUING IN IOS WAS A MAJOR CAUSE OF I/O DELAY

Finding: Queuing in the I/O Supervisor (IOSQ) was a major cause of the I/O delay with the volume.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the device.

Logic flow: The following rules cause this rule to be invoked:
DAS200: Volume with the worst overall performance

Discussion: IOSQ time is the time from the issuance of a STARTIO macro until the Start SubChannel (SSCH) instruction is issued. After the STARTIO macro is issued, the software determines whether the device is busy with the system on which the STARTIO macro was issued (that is, whether there is an available Unit Control Block (UCB) for the device). If the device is not busy with this system, the SSCH instruction is issued. However, if the device is busy with this system (a UCB is available), the I/O request is queued. Thus, IOSQ time always means that the device is unable to handle additional requests from this system.

This discussion of IOSQ time does not always apply to Parallel Access Volumes (PAVs)¹. With PAV devices, MVS creates multiple UCBs for each device, depending on how many “alias devices” have been defined. The multiple UCBs allow multiple active concurrent I/Os on a given device when the I/O requests originate from the same system². Using PAVs can dramatically improve I/O performance by nearly eliminating IOSQ.

Rule DAS250 is produced if the average IOSQ time accounted for a significant percent of the device response time.

Please refer to Rule DAS150 for additional information regarding IOSQ delay.

¹PAV devices are available with Enterprise Storage Server (ESS). With PAV devices, a “base device” address is defined, and a UCB is associated with this base address. “Alias device” addresses can be defined and UCBs are associated with the alias device addresses.

²Multiple Allegiance allows multiple active concurrent I/O operations on a given device when the I/O requests originate from different systems. The Multiple Allegiance feature is available with Enterprise Storage Server (ESS).

Rule DAS260: DEVICE DISCONNECT WAS MAJOR CAUSE OF I/O DELAY

Finding: CPExpert determined that device disconnect (DISC) time was the major cause of delay in delay in DASD response to critical applications for the device.

Impact: This finding may have a MEDIUM IMPACT or HIGH IMPACT on the performance of the critical application, on the device, and on the performance of other volumes attached to the cache controller.

Logic flow: The following rules cause this rule to be invoked:
DAS200: Volume with the worst overall performance

Discussion: DISC means that there is some delay that is often (but not always) associated with a mechanical movement during which the device disconnects from the control unit (or the control unit disconnects from the channel).

With legacy systems (e.g., 3380 drives attached to 3990-2 control units), the DISC time of most concern was associated with seek (arm movement) and rotational position sensing (time waiting for the disk platter to rotate to the location where desired data resides). Considerable performance improvement efforts were directed at reducing the seek activity and reducing the rotational position sensing (RPS)¹ delays for the legacy systems. These two mechanical delays still exist for most modern *redundant array of independent disks* (RAID)² systems, but their impact can not be directly reduced with normal methods.

With modern disks, data is cached into device cache buffers that contain data read from a track on the disk platter. Using device cache buffers containing the track data eliminated the multiple-RPS delays caused by a path busy when the device tried to reconnect. Required data is read into the device cache buffer during a single rotation and stored until a path is available to transfer the data.

In addition to the cache buffer design, modern control units such as the 3990-6 or 2105 have very large cache memory installed. With cache in the

¹RPS delays were caused by a path not being available when the required data came under a device read head. Since a path was not available, the device could not reconnect to the channel or control unit. Consequently, data could not be read and transmitted, and another rotation of the platter was experienced until the data again came under the device read head. Multiple rotations might be required, depending on the busy level of the path.

²An array is an ordered collection of physical devices (disk drive modules) that are used to define logical volumes or devices.

control units, data to be read can be transferred in a variety of ways, depending on where the data resides.

For a read operation, desired data often is found in the control unit cache. If the required data is in cache, the data can be transferred between the control unit cache and the channel, and this transfer is done at channel speed. If the required data is not in cache, the data can be transferred between the device and channel (and concurrently placed into the control unit cache for subsequent access).

For write operations, data can be placed into Non-volatile Storage (NVS) as a part of the control unit. Write operations normally end as the data to be written is placed in the NVS; and the storage processor writes the data to the device asynchronous with other activity (as a “back end” staging operation). See subsequent discussion for more detail about read and write operations.

The storage director can simultaneously transfer data between the channel and device and manage the data transfer of different tracks between the cache and channel, and the cache and the device. With large amounts of cache memory, a high percent of data accesses normally will be resolved from the fast cache memory and the relatively slow device will not cause significant delays.

As a result of the above improvements, DISC time for modern systems is a result of *cache read miss* for read operations, back-end staging delay for write operations, peer-to-peer remote copy (PPRC) operations, and other miscellaneous reasons³. DISC time often can be very small with adequate cache. For example, there would be zero disconnect time for a cache read hit (the record was found in the cache). However, DISC time can be large and can cause serious delay to I/O operations.

Suggestion: Please refer to Rule DAS160 for further information about DISC time.

³ Artis has described a “sibling PEND” condition that results from collisions within the physical disk subsystem of RAID devices. See “Sibling PEND: Like a Wheel within a Wheel,” www.cmg.org/cmgpap/int449.pdf. While this condition is titled “sibling PEND,” the time properly belongs in DISC time, rather than PEND time .

Rule DAS270: THERE DID NOT APPEAR TO BE A PROBLEM WITH VOLUME PERFORMANCE

Finding: The performance of the volume having the "worst" overall performance was not poor in any area evaluated by CPEXpert.

Impact: This finding has little impact, other than to indicate that there may be no problem with the device.

Logic flow: The following rules cause this rule to be invoked:
 DAS200: Volume with worst overall performance

Discussion: CPEXpert checked seek time, missed RPS reconnect time, PEND time, CONN time, missed cache hits, and IOSQ time. None of these areas accounted for the delay for the volume selected as the volume with the "worst" overall performance.

Suggestion: This condition most likely would occur when the guidance provided to CPEXpert is too restrictive. For example, if a large number of volumes are excluded from analysis (using the EXCLUDE option), the remaining volumes may have no particular problem. However, the logic is designed to select a "worst" volume, regardless of whether that volume actually has problems.

Rule DAS280: APPLICATIONS ACCESSING VOLUME WITH WORST PERFORMANCE

Finding: CPExpert identifies the applications (other than the "loved one" workload accessing the volume with the worst performance.

Impact: This finding is used to select applications for potential rescheduling, file movement, or design changes. This rule is applicable only if CPExpert is performing "expanded" analysis of the data.

Logic flow: The following rules cause this rule to be invoked:
DAS200: Volume with worst overall performance

Discussion: CPExpert lists the applications (other than the "loved one" workload) accessing the volume with the worst performance during the entire measurement period being analyzed. The list is ordered descendingly by I/O operations, so you can identify the applications that probably had the most impact on the volume.

As described in Section 5, the SMF Type 30(Interval) information might not be synchronized with the SMF Type 70(series) information. In this case, the identification of applications may not correctly identify the applications with the most I/O operations to the volume. However, if the analysis produces the same results after analyzing more than one day's measurement data, you can be more comfortable that the applications are correctly identified.

In some cases, the Type 30(Interval) data does not reflect all I/O activity to a device. This normally is caused by the Interval records not being written for all address spaces (such as started tasks running in SYSSTC service class). Consequently, you should use the output from DAS180 with care.

However, if the analysis produces the same results after analyzing more than one day's measurement data, you can be more comfortable that the applications are correctly identified.

Suggestion: In many cases, once the applications are identified, either applications personnel or systems personnel will verify whether the particular applications are responsible for the majority of the I/O operations to the particular volume.

Once you are comfortable that the applications are responsible for the majority of the I/O operations to the "worst" performing device, you can take action to (1) reschedule the applications to minimize contention, (2) examine the files referenced by the applications, or (3) change the way in which the applications access their files.

**Rule DAS285: APPLICATIONS ACCESSING VOLUME DURING
MEASUREMENT INTERVAL WITH WORST PERFORMANCE**

Finding: CPEXpert identifies the applications (other than the "loved one" workload) accessing the volume during the the measurement interval with the worst performance.

Impact: This finding is used to select applications for potential rescheduling, file movement, or design changes.

Logic flow: The following rules cause this rule to be invoked:
DAS200: Volume with worst overall performance

Discussion: CPEXpert lists the applications accessing the volume with the worst performance during the measurement period with the worst performance being analyzed. The list is ordered descendingly by I/O operations, so you can identify the applications that probably had the most impact on the volume.

As described in Section 5, the SMF Type 30(Interval) information might not be synchronized with the SMF Type 70(series) information. In this case, the identification of applications may not correctly identify the applications with the most I/O operations to the volume. However, if the analysis produces the same results after analyzing more than one day's measurement data, you can be more comfortable that the applications are correctly identified.

In some cases, the Type 30(Interval) data does not reflect all I/O activity to a device. This normally is caused by the Interval records not being written for all address spaces (such as started tasks running in SYSSTC service class). Consequently, you should use the output from DAS180 with care.

However, if the analysis produces the same results after analyzing more than one day's measurement data, you can be more comfortable that the applications are correctly identified.

Suggestion: In many cases, once the applications are identified, either applications personnel or systems personnel will verify whether the particular applications are responsible for the majority of the I/O operations to the particular volume.

Once you are comfortable that the applications are responsible for the majority of the I/O operations to the "worst" performing device, you can take

action to (1) reschedule the applications to minimize contention, (2) examine the files referenced by the applications, or (3) change the way in which the applications access their files.

Rule DAS287: OTHER APPLICATIONS DID NOT REFERENCE THE VOLUME

Finding: CPExpert has determined that only the "loved one" workload referenced the volume during the interval when the device had the worst performance.

Impact: This finding has little impact, other than to indicate that the device performance is a function of the "loved one" workload, rather than performance problems being caused by other workload.

Logic flow: The following rules cause this rule to be invoked:
DAS200: Volume with worst overall performance

Discussion: CPExpert examined the workload accessing the device during the period of "worst" performance. Only the "loved one" workload accessed the device during this measurement interval.

As described in Section 5, the SMF Type 30(Interval) information might not be synchronized with the SMF Type 70(series) information. In this case, the identification of applications may not correctly identify the applications with the most I/O operations to the volume. However, if the analysis produces the same results after analyzing more than one day's measurement data, you can be more comfortable that the applications are correctly identified.

In some cases, the Type 30(Interval) data does not reflect all I/O activity to a device. This normally is caused by the Interval records not being written for all address spaces (such as started tasks running in SYSSTC service class). Consequently, you should use the output from DAS180 with care.

However, if the analysis produces the same results after analyzing more than one day's measurement data, you can be more comfortable that other workload did not reference the device.

Suggestion: This finding implies that performance problems for the device are not caused by contending workloads. Performance tuning actions must be directed toward the "loved one" workload.

Rule DAS290: PERFORMANCE CHARACTERISTICS OF SIGNIFICANT VOLUMES

Finding: CPExpert identifies the performance characteristics of the most significant volumes accessed by the "loved one" workload (including the volume selected as the "worst" performing volume).

Impact: This finding is used to assess the importance of the "worst" performing device from the perspective of the "loved one" workload, and to determine whether other devices offer significant performance improvement potential.

Logic flow: This is a basic finding. There are no predecessor rules.

Discussion: CPExpert lists basic characteristics of the volumes having the most potential for improvement. The list includes the volume selected as the "worst" performing device, so that you can appreciate the relative performance improvement potential between the "worst" volume and other volumes on the list.

Volumes are selected as having improvement potential only if their response time exceeds the average for their device type. This "screening" is done for each measurement interval being analyzed. The data presented by Rule DAS190 reflects the average per-second delays only during measurement intervals when the device I/O performance was worse than the average for its device type.

The "weighted delays" value is a relative measure of the performance improvement potential of the volume. The absolute values in the column are not particularly meaningful. Rather, the values should be compared to each other to assess the relative performance impact of each volume.

For example, the "worst" volume might have a "weighted delays" value of 1000. If the "next worst" device had a "weighted delays" value of 950, you may wish to direct CPExpert to examine the next worst device in more detail (this would be accomplished by using DASGUIDE to "exclude" the worst volume).

On the other hand, suppose the "next worst" device had a "weighted delays" value of only 400. You probably would not wish to examine the next worst in more detail, since little performance benefit would be gained from tuning actions directed to the next worst volume.

It is possible that a volume may have a significant improvement potential in a particular measurement interval, but not be the volume with the most overall potential for improvement. This situation can arise because the analysis is directed toward the volumes with the **most overall** performance improvement potential.

Suggestion: You should use the information displayed by Rule DAS190 to assess the relative importance of the "worst" performing device compared with the performance improvement potential of the other devices.

Rule DAS300: Perhaps shared DASD caused performance problems

Finding: CPExpert believes that accessing conflicts caused by sharing DASD between systems or MVS images may have caused performance problems.

Impact: This finding is used to assess whether sharing DASD between systems or MVS images caused performance problems.

Logic flow: The following rules cause this rule to be invoked:

- DAS100: Volume with the worst overall performance
- DAS110: Seeking was major cause of I/O response delay
- DAS120: Missed RPS reconnect was major cause of I/O delay
- DAS130: Large PEND time was major cause of I/O delay
- DAS150: Missed cache read hits was major cause of I/O delay

Discussion: DASD volumes can be shared between systems or between MVS images. Sharing of the DASD volumes might be implemented to allow backup of data, to facilitate recovery or restart, to permit transfer of data from one system to another, etc.

In some situations, sharing DASD volumes has little impact on performance (for example, few I/O operations might be directed to the shared volumes from potentially conflicting systems).

In other situations, sharing DASD volumes can have a significant impact on the performance of the shared volumes, and consequently, on the performance of the applications accessing the shared volumes.

CPExpert can perform an analysis of conflicts between DASD shared between systems or MVS images. The analysis performed by CPExpert is not intended to identify an isolated performance problem. Rather, CPExpert attempts to identify those problems that **continually** cause shared DASD performance problems. Shared DASD analysis is invoked by specifying **%LET SHARED = Y;** in USOURCE(DASGUIDE). Shared DASD analysis is an option, because more processing is required to perform shared DASD analysis.

CPExpert performs the following processing if you have indicated that an analysis of potential conflicts between shared DASD should be performed:

- CPExpert determines whether the "worst" devices selected for detailed analysis are shared with another system. If so, CPExpert performs an analysis of potential conflicts caused by shared DASD.

-
- CPExpert identifies other systems that reference the "worst" device. This identification is accomplished by analyzing the SMF Type 74 data in the performance data base relating to all other systems. The SMF Type 74 data contain the VOLSER for each device referenced. CPExpert simply selects SMF Type 74 information for the systems that reference the VOLSER of the "worst" device. This information is retained for more detailed analysis about potential conflicts.

There is a potential (but very unlikely) problem with this method of identifying devices shared between systems. Multiple systems in the performance data base could use the same VOLSER to identify different devices. This could happen if the devices were not shared between systems.

For example, suppose that CPExpert had identified PAGE01 as the "worst" device. Several system in the performance data base could reference VOLSER PAGE01, but the devices with VOLSER PAGE01 could be unique to each system. CPExpert would assume that all references by other systems to PAGE01 applied to the "worst" device being analyzed. The references could apply to a totally different device, and the other systems might not even share DASD with the system being analyzed.

If this should be a problem (that is, if the DASD Component reports shared DASD conflicts with systems that do not share the device being analyzed), simply ignore the analysis produced by CPExpert¹.

- Once CPExpert has identified all systems that reference the "worst" device, CPExpert analyzes the DASD I/O characteristics of these systems with respect to the "worst" device. As described earlier, the analysis makes a basic assumption that the I/O activity from the different systems is random among the systems (for example, the code assumes that the I/O activity of System B is independent from the I/O activity of System A).

CPExpert will produce Rule DAS300 to list statistics relating to potential conflicts, by system, by volume, and by RMF measurement interval. The following example shows sample output from Rule DAS300:

¹We do not feel that this problem will be common. It is described only to alert you to a potential incorrect analysis. If any user encounters this problem and it becomes annoying, code can be implemented to allow users to identify specific systems that share DASD with the system being analyzed. At present, this option seems to add unnecessary complexity to the user options.

RULE DAS300: PERHAPS SHARED DASD CONFLICTS CAUSED PERFORMANCE PROBLEMS

Accessing conflicts caused by sharing VOLSER PPVOL1 between systems might have caused performance problems for the device during the measurement intervals shown below. Conflicting systems had the indicated I/O rate, average CONN time per second, average DISC time per second, average PEND time per second, and average RESERVE time to the device. Even moderate CONN, DISC, or RESERVE can cause delays to shared devices.

MEASUREMENT INTERVAL	I/O RATE	MAJOR PROBLEM	OTHER SYSTEM	-----OTHER SYSTEM DATA-----	I/O RATE	CONN	DISC	PEND	RESV
8:30- 8:45,22OCT2001	131.6	PEND TIME	J80		37.3	0.042	0.003	0.065	0.000
			JF0		147.0	0.129	0.005	0.158	0.000
			JH0		368.8	0.372	0.036	0.799	0.000
			Z0		459.7	0.406	0.017	0.765	0.000
			TOTAL		1012.8	0.949	0.061	1.786	0.001
8:45- 9:00,22OCT2001	108.5	PEND TIME	J80		41.2	0.046	0.003	0.066	0.000
			JF0		195.1	0.169	0.006	0.226	0.000
			JH0		411.7	0.411	0.032	0.718	0.001
			Z0		498.9	0.432	0.015	0.795	0.001
			TOTAL		1147.0	1.058	0.056	1.805	0.002

Rule DAS300 shows, for each RMF measurement interval, the I/O rate and the major problem during the RMF interval, of the device being analyzed. The remaining data shows relevant information (I/O rate, CONN time, DISC time, PEND time, and RESERVE time) for the other systems that reference the device.

CPEXpert summarizes the other system data, into a TOTAL row for each RMF interval. In some instances, the TOTAL per second time in a particular interval will be more than one second. In the case of DISC time or PEND time, this situation is caused by multiple I/O operations being delayed for the device. This commonly happens only with Parallel Access Volume (PAV) devices.

However, the above example shows that the CONN time is larger than one second per second! CONN time is normally thought to involve data transfer between the device and the host system, and this concept is convenient for most analysis. Clearly, a device cannot be active transferring data for more than one second per second.

The CONN time actually is a hardware construct that is measured at the channel subsystem level, and includes all hardware protocol between the channel, the director port, and the control unit. The hardware protocol is a very small amount of time. However, if there are many I/O operations, and the hardware protocol must take place for each I/O operation (and takes place at several points in the I/O operation), the total hardware protocol can become more significant.

In the above example, there was a total of 1147.0 I/O operations shown in the last total line (which is the total for “other systems” referencing the device). If the hardware protocol connect time were multiplied by 1147 I/O operations per second, it is easy to appreciate that the total hardware protocol connect time from the multiple systems would add appreciably to the connect time.

Thus, while the connect time for a particular device cannot exceed one second per second, once the hardware protocol connect time is added to the device connect time, the total can exceed one second connect time per second for a very active device.

Suggestion: You should use the information displayed by Rule DAS300 to assess the significance of the performance problems caused by shared DASD.

Rule DAS310: Seeking probably was caused by shared DASD conflicts

Finding: CPExpert believes that the seeking performance problems were caused by shared DASD conflicts.

Impact: This finding is used to assess whether sharing DASD between systems or MVS images caused performance problems.

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance
DAS110: Seeking was the major cause of I/O response delay
DAS300: Shared DASD caused performance problems

Discussion: If CPExpert determines that seeking was the major cause of I/O response delay and if the device is shared, CPExpert analyzes other systems in the performance data base which share the volume.

Seek delays can occur if the arm of the device has been moved by System B when System A attempts to access a cylinder. In this case, the DASD I/O operation from System A must move the arm to the desired cylinder, as a SEEK operation.

CPExpert computes the time required to perform SEEKS on System A. If the computed SEEK time is a major cause of performance problems, CPExpert analyzes the data from System B to determine whether System B generates a large number of I/O operations to the device.

If System B does **not** generate a relatively large number of I/O operations to the device, CPExpert concludes that there is **not** a conflict. If System B **does** generate a relatively large number of I/O operations to the device, CPExpert concludes that there **is** a conflict caused by sharing the device.

- There is little doubt about the validity of the first conclusion: if System B does not use the device, System B clearly cannot cause seek problems for System A.
- The second conclusion is based on an assumption: a large number of I/O operations from System B to the shared device will cause the seek problems for System A.

To be absolutely correct, CPExpert should process the configuration definitions for System B, process System B's channel and device information, and compute seek information for System B. CPExpert could

then determine whether System B also experienced a high seek rate for the device. If both System A and System B experienced a high seek rate, CPEXpert could be absolutely sure that there was a shared DASD conflict. This approach would unnecessarily use system resources and would be cumbersome to implement.

Consequently, CPEXpert makes the assumption that I/O operations to the device are random between System A and System B. CPEXpert thus can conclude that if System A experiences a high seek rate and System B significantly uses the device (exhibited by a high I/O rate), then System B must also experience a high seek rate. To assume otherwise would require that I/O from System B be coordinated with the I/O from System A, such that System B does not experience seeking similar to System A.

Since both System A experiences a high seek rate and System B experience a high I/O rate, CPEXpert concludes that there is a conflict caused by the shared DASD.

Suggestion: You should use the information displayed by Rule DAS300 to assess the significance of the performance problems caused by shared DASD. Please refer to the suggestions associated with Rule DAS300 for alternative actions you may consider.

Rule DAS330: Large PEND time probably was caused by shared DASD conflicts

Finding: CPEXpert believes that the large PEND time performance problems were caused by shared DASD conflicts.

Impact: This finding is used to assess whether sharing DASD between systems or MVS images caused performance problems.

Logic flow: The following rules cause this rule to be invoked:
 DAS100: Volume with the worst overall performance
 DAS130: Major cause of I/O delay was PEND time
 DAS132: Large PEND time may be caused by other system
 DAS300: Shared DASD conflicts caused performance problems

Discussion: If CPEXpert determines that large PEND time was the major cause of I/O response delay and if the device is shared, CPEXpert analyzes other systems in the performance data base which share the volume.

The PEND time is a measure of the time an I/O operation waited (1) for a channel path to become available, (2) for the controller to become available, or (3) for the device because the device was busy to another system.

- If the I/O operation waited for a channel path, the I/O delay was caused by activity on the system issuing the I/O operation. Since channel paths are not shared between systems, the delay must occur at the controller or device level if the delay was caused by sharing DASD.
- If the I/O operation waited for the controller to become available, the I/O delay could be caused by activity on the system issuing the I/O operation or could be caused by another system which was using the controller.
- If the I/O operation waited for the device, the I/O delay was caused by another system which was using the device. The I/O Supervisor (IOS) on the system issuing the I/O operation could not know that the device was busy with another system (the device Unit Control Block would not reflect the "busy" status). Consequently, the IOS would believe the device was available and would issue the I/O operation. The I/O operation would wait while the device was busy to the other system. This time would be reflected as PEND time.

CPEXpert analyzes the I/O activity from other systems sharing the device. The I/O activity measures used by CPEXpert consists of the DISC time and CONN time from other systems. CPEXpert substitutes the total RESERVE time for other systems if the total RESERVE time for other systems is greater than the I/O activity time.

CPEXpert concludes the high PEND time is caused by other systems sharing the device if the device I/O activity (or total RESERVE time) is more than 25% of the PEND time experienced by the system being analyzed.

Suggestion: You should use the information displayed by Rule DAS300 to assess the significance of the performance problems caused by shared DASD. Please refer to the suggestions associated with Rule DAS300 for alternative actions you may consider.

Rule DAS360: Missed cache hits probably were caused by shared DASD conflicts

Finding: CPExpert believes that the missed cache hits performance problems were caused by shared DASD conflicts.

Impact: This finding is used to assess whether sharing DASD between systems or MVS images caused performance problems.

Logic flow: The following rules cause this rule to be invoked :

- DAS100: Volume with the worst overall performance
- DAS160: Missed cache read hits was major cause of I/O delay
- DAS300: Shared DASD conflicts caused performance problems

Discussion: If CPExpert determines that cache missed read hits was the major cause of I/O response delay and if the device is shared, CPExpert analyzes other systems in the performance data base which share the volume.

Cache read hits can be caused by another system if the other system reads data from devices attached to the cached controller. If the data required by the other system is not in cache, the data must be brought into the cache to replace data already in the cache. The data being replaced might subsequently be required by the system being analyzed, and thus missed cache read hits would occur.

- If the data required by the other system is already in the cache, the data transfer will be reflected as CONN time by the other system.
- If the data must be brought into the cache, the mechanical time for the device will be reflected in DISC time. This is the time to position the arm (seek time) and to rotate the disk to the proper sector (latency). Rule DAS110 and Rule DAS310 would be produced if CPExpert estimated that seeking was a problem. Consequently, the residual time examined by Rule DAS360 is the latency delay.

In either of the above situations with cached shared devices, CPExpert analyzes the amount of I/O operations from System A and System B to the cached controller.

- If System B does **not** generate a relatively large number of I/O operations to the device, CPExpert concludes that there is **not** a conflict.

-
- If System B **does** generate a relatively large number of I/O operations to the device, CPEXpert concludes that there **is** a conflict caused by sharing the device. The rationale for these conclusions is the same as was discussed in Rule DAS310 (seeking probably was caused by shared DASD conflicts).

The standard analysis performed by CPEXpert may not detect a cache problem under two possible scenarios.

- If only one cached controller is attached to System A, CPEXpert may not detect a problem with the device. This is because the logic employed by CPEXpert selects devices with the most performance improvement within each **type of device** and then selects the **overall** "worst" devices for detailed analysis.

CPEXpert considers cached devices to be a unique type. If all devices on the cached controller received bad service caused by shared cached problems, CPEXpert may not detect a performance problem with the cached devices. This is because all devices in the "device type" could have roughly equal poor service and no device would be **significantly** worse than the other devices in the device type. Consequently, CPEXpert might not select any of the cached devices for detailed analysis

If there are multiple cached controllers, there would be a larger number of "candidate" devices, and the standard analysis performed by CPEXpert is more likely to identify any problem caused by shared cache controllers.

- The analysis performed by CPEXpert may not detect a cache problem if the cache is being replaced with data from another volume. It is possible that System B could cause data from another volume to be loaded into the controller's cache. This volume could be a volume not be flagged as the "worst" performing volume when analyzing performance from the perspective of System A. From System A's perspective, accesses to the worst volume would simply not find required data in cache.

Without analyzing the IOCP information for all systems, CPEXpert cannot determine which volumes are attached to which controllers. Consequently, CPEXpert cannot at present relate (1) poor performance for one volume on System A and (2) I/O operations to a different volume by System B.

If you suspect problems because of shared cached devices, you can direct CPEXpert to analyze the specific devices (using the SELECT option in DASGUIDE). CPEXpert will then analyze only the devices selected.

Suggestion: You should use the information displayed by Rule DAS300 to assess the significance of the performance problems caused by shared DASD. Please refer to the suggestions associated with Rule DAS300 for alternative actions you may consider.

Rule DAS380: Applications potentially causing shared DASD conflicts

Finding: CPExpert lists the applications that might cause shared DASD conflicts. |

Impact: This finding is used to assess whether sharing DASD between systems or MVS images caused performance problems, and to identify specific applications that might cause conflicts. |

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance
DAS300: Shared DASD caused performance problems

This finding is produced only if the CPExpert modification to MXG or MICS is installed, so that SMF Type 30(DD) information is available.

Discussion: If CPExpert determines that there may be DASD conflicts caused by sharing DASD between systems of MVS images, and if the CPExpert modification to MXG or MICS is installed to collect application/device data, CPExpert performs further analysis.

CPExpert analyzes the SMF Type 74 information to acquire the device numbers associated with each VOLSER on every system. This processing is required because, while all systems would refer to a particular device by the VOLSER of the device, each system could reference the device using a different device number.

CPExpert creates the DASD30DD data set from the Type 30 records, using the modification to MXG or MICS. VOLSER information is not retained in the SMF Type 30 records; device numbers are retained in the SMF Type 30 records. Consequently, CPExpert must associate the unique device number used by each system, to the VOLSER of the device.

CPExpert then analyzes the DASD30DD data set created by the modification to MXG or MICS. CPExpert selects all job steps executing on **other** systems that referenced the device number of the device with the worst performance. For example, if CPExpert is analyzing System A, CPExpert would select all job steps executing on System B, System C, etc. if those job steps referenced the device with the worst performance. |

CPExpert includes only references for time intervals when earlier analysis had concluded that shared DASD conflicts might cause a performance problem.

CPEXpert produces a list of applications (listed by SMF Type 30 interval start and end times), ranked descendingly by I/O operations to the device with the worst performance.

As described in Section 5, the SMF Type 30(Interval) information might not be synchronized with the SMF Type 70(series) information. In this case, the identification of applications may not correctly identify the applications with the most I/O operations to the volume. However, if the analysis produces the same results after analyzing more than one day's measurement data, you can be more comfortable that the applications are correctly identified.

In some cases, the Type 30(Interval) data does not reflect all I/O activity to a device. This normally is caused by the Interval records not being written for all address spaces (such as started tasks running in SYSSTC service class). Consequently, you should use the output from DAS180 with care.

Suggestion: Rule DAS380 is provided so you can identify the applications potentially causing shared DASD conflicts. You may wish to take the following actions:

- Reschedule the applications that cause shared DASD conflicts to a different processing time. The shared DASD conflicts may be eliminated or reduced if the applications can be rescheduled.
- Review the data sets on the volume to determine whether the data sets can be moved to other shared volumes (or to non-shared volumes). The shared DASD conflicts may be reduced or eliminated if the **data sets** are not required by systems sharing the volumes.
- Modify the applications causing shared DASD conflicts to reduce or eliminate their use of the data on the shared volume.

Rule DAS385: Applications potentially causing worst shared DASD conflicts

Finding: CPExpert lists the applications that might cause shared DASD conflicts during the RMF measurement intervals with the worst performance. |

Impact: This finding is used to assess whether sharing DASD between systems or MVS images caused performance problems, and to identify specific applications that might cause the worst conflicts. |

Logic flow: The following rules cause this rule to be invoked:
DAS100: Volume with the worst overall performance
DAS300: Shared DASD conflicts caused performance problems

This finding is produced only if the CPExpert modification to MXG or MICS is installed, so that SMF Type 30(DD) information is available.

Discussion: If CPExpert determines that there may be DASD conflicts caused by sharing DASD between systems of MVS images, and if the CPExpert modification to MXG or MICS is installed to collect application/device data, CPExpert performs further analysis. This analysis is described in Rule DAS380.

CPExpert further processes the information selected from the DASD30DD data set, to select job steps that referenced the device with the worst performance, **during the RMF measurement interval in which the device had the worst performance**. CPExpert includes only references for time intervals when earlier analysis had concluded that shared DASD conflicts might cause a performance problem. |

CPExpert produces a list of applications (listed by SMF Type 30 interval start and end times), ranked descendingly by the access to the device with the worst performance. The SMF Type 30 interval data are prorated, if necessary, to reflect the approximate number of I/O operations during the worst RMF interval. Please refer to Section 5 (Chapter 2.1) for a discussion of prorating issues. |

As described in Section 5, the SMF Type 30(Interval) information might not be synchronized with the SMF Type 70(series) information. In this case, the identification of applications may not correctly identify the applications with the most I/O operations to the volume. However, if the analysis produces the same results after analyzing more than one day's measurement data, you can be more comfortable that the applications are correctly identified. |

In some cases, the Type 30(Interval) data does not reflect all I/O activity to a device. This normally is caused by the Interval records not being written for all address spaces (such as started tasks running in SYSSTC service class). Consequently, you should use the output from DAS180 with care.

Suggestion: Rule DAS385 is provided so you can identify the applications potentially causing the worst shared DASD conflicts. You may wish to take the following actions:

- Reschedule the applications that cause shared DASD conflicts to a different processing time. The shared DASD conflicts may be eliminated or reduced if the applications can be rescheduled.
- Review the data sets on the volume to determine whether the data sets can be moved to other shared volumes (or to non-shared volumes). The shared DASD conflicts may be reduced or eliminated if the **data sets** are not required by systems sharing the volumes.
- Modify the applications causing shared DASD conflicts to reduce or eliminate their use of the data on the shared volume.

Rule DAS390: Shared DASD conflicts did not cause performance problems

Finding: CPExpert did not detect performance problems that were caused by shared DASD conflicts.

Impact: This finding is used to assess whether sharing DASD between systems or MVS images caused performance problems.

Logic flow: This is a basic finding. There are no predecessor rules.

Discussion: Parameters in SMF data indicated that the device with the worst performance was shared between systems or MVS images.

CPExpert analyzed data from all systems in the performance data base. The analysis did not indicate that the performance problems were caused by sharing the device.

This finding could result because the device was not configured on-line to another system, even though the SMF data indicated that the device was shared.

Alternatively, the other system might not have generated a large amount of I/O activity to the device.

In either case, CPExpert does not believe that sharing DASD caused conflicts.

Suggestion: There are no suggestions with the finding. Rule DAS390 is provided simply for information purposes.

Rule DAS400: Access characteristics of significant data sets

Finding: CPExpert identifies the access characteristics of data sets managed by DFSMS (and reported in SMF Type 42 records), for the volumes that have the most potential for performance improvement.

Impact: This finding is used to assess the data sets that reside on the volumes with significant performance improvement potential.

Logic flow: This is a basic finding. There are no predecessor rules.

Discussion: If SMF Type 42 (Data Set Statistics) information is available¹ in a MXG performance data base, the DASD Component will process the MXG TYPE42DS file. The DASD Component select information describing data sets that were referenced during RMF measurement intervals in which the poorly performing device(s) exceeded the average performance. The result from this analysis is displayed in Rule DAS400.

There are several considerations with respect to the information produced by Rule DAS400:

- There can be many data sets referenced on the devices (hundreds or even thousands of data sets can be referenced). It is not helpful to have a large number of data sets listed. Consequently, CPExpert provides two guidance variables to control how many data sets are listed by Rule DAS400:
- The **LIST42DS** variable can be used to limit the number of data sets listed in any RMF measurement interval, for a particular device. Only the number of data sets specified by the LIST42DS variable will be listed individually, and information about any remaining data sets will be summarized and listed on a single line.

The default value for the LIST42DS variable is 10, indicating that 10 data sets will be listed for each RMF interval in which a poorly-performing device had a performance problem. You can alter this number of data sets listed by using the LIST42DS guidance variable.

- Some data sets might have a low I/O activity and would not be interesting to analyze. CPExpert provides the **MIN42PCT** variable to limit the number of data sets listed in any RMF measurement interval,

¹%LET TYPE42DS = Y; was specified in USOURCE(GENGUIDE).

for a particular device. Only data sets having a percent of activity for the total volume greater than the percent specified by the MIN42PCT variable will be listed individually, and any remaining data sets will be summarized and listed on a single line.

The default value for the MIN42PCT variable is 0.1, indicating that data sets will not be listed individually for any RMF interval unless the data set intensity of access (I/O rate * response time) was greater than 0.1% of the total volume intensity of access. You can alter this percent of data sets listed by using the MIN42PCT guidance variable.

For example, if you specified %LET MIN42PCT = 25, a maximum of 4 data sets would be listed in any RMF interval since no more than 4 data sets could have 25% or higher access intensity.

Please note that, regardless of the data set access intensity of any particular data set, only the number of data sets specified by the **LIST42DS** guidance variable (described earlier) will be listed. This means that there are two ways to limit the number of data sets listed: (1) the LIST42DS which limits the number of data sets listed in any RMF interval, and (2) the **MIN42PCT** guidance variable which limits the data sets listed to those that exceed the specified percent.

- Some devices might have data sets that are not managed by DFSMS. In this case, information would not be available in SMF Type 42 records and the data sets could not be reported by CPExpert.
- Depending on the timing in which SMF writes Type 42 records, some information might not be represented in the RMF Type 74 intervals by which the DASD Component analyzes DASD performance problems. Consequently, the information presented in Rule DAS400 might not correspond to the information presented in other rules (for example, might not correspond to information presented in Rule DAS100).
- SMF writes Type 42 (Data Set Statistics) records when (1) a DASD data set is closed, or (2) immediately after the recording of the SMF Type 30 interval record. There are two implications of this SMF write characteristic:
 - A data set likely would be closed at some time between SMF writing Type 30 interval records. Consequently, some Type 42 (Data Set Statistics) records would not correspond to the SMF Type 30 interval, and the data set information reported by the DASD Component might not correspond to the information reported for the device by other rules.
 - The SMF Type 30 recording interval controls when Type 42 (Data Set Statistics) records are written (for those data sets that are OPEN when

the Type 30 record is written). However, the SMF Type 30 recording interval might not correspond with the RMF recording interval. If the SMF and RMF writing is synchronized, the SMF Type 30 recording interval would be synchronized with the RMF recording interval. However, the synchronization option might not have been selected.

Even if the synchronization option had been selected, a different recording interval frequency could have been selected for Type 30 and RMF recording. For example, RMF recording could occur at 15 minute intervals, but SMF Type 30 recording could occur at 30 minute intervals. While the intervals could be synchronized, there would be twice as many RMF recording intervals as there were SMF Type 30 (and Type 42) recording intervals.

The SMF and RMF recording considerations might cause consternation if a user were not aware of the potential problems. If the SMF and RMF recording intervals are synchronized and are for the same duration, and if data sets are OPEN for the duration of the recording interval, information presented in Rule DAS400 will correspond well with device-related information (and application-related information) presented by other rules. If these conditions are not present, the information presented by Rule DAS400 might not correspond well with information presented elsewhere by the DASD Component.

Please send a note to Don_Deese@cpexpert.com if these problems become troublesome. Perhaps the logic can be improved (by prorating the DASD42DS information, for example).

Suggestion: You should use the information displayed by Rule DAS400 to assess the data access characteristics for DFSMS-managed data sets residing on the device(s) with the most potential for performance improvement.

Rule DAS600: Excessive Control Area (CA) splits occurred

Finding: CPExpert determined that an excessive number of Control Area (CA) splits occurred for the VSAM data sets listed.

Impact: This finding is used to assess problems or potential problems with VSAM Control Area splits for the data sets listed.

Discussion: A VSAM file structure consists of one or more *Control Intervals (CIs)* and one or more *Control Areas (CAs)*.

- A **Control Interval** is a continuous area of direct access storage that VSAM uses to store logical records. When a logical record is read from direct access storage, the entire Control Interval containing the record is read into a VSAM buffer in virtual storage. The desired logical record is then transferred from the VSAM buffer to a user-defined buffer or work area. While logical records within a Control Interval may vary in length, all Control Intervals in a specific VSAM data set are of the same length.

In addition to the logical records, a Control Interval consists of free space, and control information. The free space initially is unused, and is used to accommodate inserted logical records or for changes in the length of logical records. The control information describes the amount and location of free space, and describes the length of records and how many adjacent records are of the same length.

- A **Control Area** contains one or more Control Intervals. The Control Intervals are grouped together into fixed-length contiguous areas of direct access storage. A VSAM data set consists of one or more Control Areas.

The minimum size of a Control Area is one track of DASD storage, and the maximum size of a Control Area is one cylinder. The size of a Control Area is not specified by a user; the size of the CA is *calculated* by VSAM based on the amount of space allocated to a VSAM data set.

Logical records can be inserted into a VSAM keyed sequenced data set (KSDS) or a variable relative record data set (VRRDS), or a record length can be increased. When either of the actions occur, space must be available in a Control Interval to accommodate the new record (or accommodate the increased length). The space for the new record (or

increased length) normally is obtained from the *free space* that was allocated to the Control Interval¹ when the data set was loaded.

If enough free space is not available in the Control Interval, a *Control Interval Split* occurs. When a Control Interval is split, approximately one-half of the logical records in the original Control Interval will be transferred to a new Control Interval, and the records are then deleted in the old Control Interval. Space for the new Control Interval is obtained from the free space that is associated with the Control Area. Thus, inserted records or longer records are inserted into Control Intervals using the free space associated with the Control Interval. If the Control Interval does not have sufficient space, a new Control Interval will be created, using free space in the Control Area² to which the original Control Interval belongs.

If enough free space is not available in the Control Area, a *Control Area Split* occurs. When a Control Area is split, approximately one-half of the Control Intervals from the Control Area are moved to the end of the data set. The records are then deleted from the old Control Area (requiring additional I/O operations). For example, if the CA is one cylinder, then approximately one-half of the records in the cylinder is moved to the end of the data set, and the records that were moved are then deleted from their location in the original cylinder. Depending on the size of the Control Area, considerable overhead³ can result from Control Area splits.

Additionally, locking may be involved during Control Area split processing, depending on the SHARE options associated with the VSAM data set. The below table shows the level at which locking occurs, and the condition that causes the lock level.

LOCK LEVEL	CONDITION
Control Interval	Adding a record or updating a record in place, without causing a Control Interval split
Control Area	Adding a record or updating a record in place, causing a Control Interval split, but not causing a Control Area split.
Data Set	Adding a record or updating a record in place, causing a Control Interval split, and causing a Control Area split.

¹The amount of free space reserved for each Control Interval is controlled by the FREESPACE parameter. The first value of the FREESPACE keyword specifies a percent of each Control Interval that is to be reserved for free space.

²The amount of free space reserved for each Control Area is controlled by the FREESPACE parameter. The second value of the FREESPACE parameter specifies a percent of each Control Area that is to be reserved for free space.

³Note that the overhead occurs only during the CA split. There is essentially no additional overhead after the CA split.

Considering the overhead involved and the locking that can occur, Control Area splits should be avoided if possible.

CPEXpert examines the SMF Type 64 information contained in MXG TYPE64 data set to identify VSAM data sets that have excessive Control Area splits.

CPEXpert sums the ACCASPLT variable (the number of CA splits since the data set was created) and the CASPLITS variable (the number of CA splits with the current OPEN of the data set). CPEXpert compares this sum with the **CASPLITS** guidance variable in USOURCE(DASGUIDE). CPEXpert produces Rule DAS600 when the total number of CA splits exceeds the value specified by the **CASPLITS** guidance variable.

The default value for the **CASPLITS** guidance variable is 10, indicating that CPEXpert should produce Rule DAS600 when a VSAM data set experienced more than 10 CA splits.

The following example illustrates the output from Rule DAS600:

RULE DAS600: EXCESSIVE CONTROL AREA SPLITS OCCURRED

VOLSER: RLS014. More than 10 Control Area (CA) splits occurred for the VSAM data sets listed below. CA splits cause considerable overhead during the split, and CA splits should be avoided.

SMF TIME STAMP	JOB NAME	VSAM DATA SET	TOTAL CA SPLITS	CA SPLITS THIS OPEN
10:30,29AUG2001	CICS2AGC	RLSADSW.VF03D.DATAENDB.DATA.....	115	9
10:30,29AUG2001	CICS2AGC	RLSADSW.VF04D.DATAENDB.DATA.....	117	11
10:30,29AUG2001	CICS2AGA	RLSADSW.VF07D.ITEMACT.DATA.....	63	49

Although not shown in this example, CPEXpert also reports the total number of inserts to the VSAM data set and inserts for the current OPEN of the VSAM data set.

Suggestion: The action taken to reduce CA splits requires a reorganization of the data set, often with different values for the FREESPACE parameter. If CA splits occur frequently for a VSAM data set, you should consider the following alternatives:

- **Reduce CI free space and increase CA free space.** Unless inserts (or increases in record length) are evenly distributed across Control Intervals, you should reduce (or eliminate) the Control Interval FREESPACE amount and increase the Control Area FREESPACE amount. Control Interval splits do not cause much overhead (little data is moved and the movement remains within the Control Area). Reducing

or eliminating the Control Interval FREESPACE amount eliminates wasted DASD space, unless inserts take place for most Control Intervals. Adding free space to the Control Area will often significantly reduce the number of Control Area splits and significantly reduce the overhead associated with the Control Area splits.

If new records will be evenly distributed throughout the data set, control area free space should equal the percentage of records to be added to the data set after the data set is loaded (specify FREESPACE (0 nn), where nn equals the percentage of records to be added.)

Unfortunately, simply increasing the amount of Control Area free space might not be a good solution. Specifying too much free space can result in more direct access storage required to contain the data set and much of this space might be wasted if inserts are clustered.

- **Allow the insert pattern to control where splits occur.** If insertions will be unevenly distributed throughout the data set and you cannot tell where the insertions are likely to occur, you can specify a small amount of free space. In this case, splits will occur *in those areas where insertions occur*. This option reduces the amount of wasted DASD space caused by a large amount of unused free space.
- **Allocate CA free space where clustering occurs.** If you can determine the clustering nature of the record insertions and the data set is large, you might consider loading the data set in stages⁴, with free space specified in areas where clustering is expected.
- **Spread the overhead.** If the VSAM data set is used **primarily by on-line applications**, you might specify a relatively small primary allocation, so that the impact of Control Area splits would be spread across the transactions encountering the Control Area splits. If a Control Area were one cylinder, the considerable overhead of splitting the Control Area would be incurred by the transaction that caused a Control Area split to occur. This might be a relatively random event (depending on the characteristics of other transactions and which user submitted the transactions), and would result in unpredictably poor response time. One way to reduce the impact on any particular on-line user would be to make the Control Area relatively small (for example, use a *primary allocation unit* of only a track, which would cause the Control Area to be only a track) so that when a CA split did occur, the immediate impact would not be as significant.

⁴ IBM shows an example of this approach in *DFSMS: Using Data Sets* (Section 2.5.3.2: Altering the Free Space Specification When Loading a Data Set)

-
- **You can not determine the clustering nature of the record insertions and the data set is relatively small.** If this is the situation, you probably should ignore the CA splits reported by DAS600, since the CA splits should not continue to occur. If the CA splits **should** continue to occur, you probably will see the data set allocation begin to acquire additional extents. This is because the size of the CA is determined by VSAM based on the primary and secondary allocation. In this situation, CPExpert will produce Rule DAS604 (Excessive secondary extents were allocated).
 - If none of the above actions are appropriate, you can change the **CASPLITS** guidance variable in USOURCE(DASGUIDE). Section 3 describes how to change the CASPLITS guidance variable if you feel that Rule DAS600 is produced too often, or if you do not wish to take action when only 10 CA splits occur.
 - Alternatively, you can exclude the reported VSAM data sets from analysis. Section 3 describes how to exclude VSAM data sets from analysis. However, you should be aware that no analysis of potential VSAM problems will be performed on data sets that are excluded from analysis.

Reference: *DFSMS: Using Data Sets* (SC26-7339 for OS/390; SC26-7410 for z/OS)
Section 2.11.4.2.3: Control Interval Splits
Section 2.5.2: Optimizing Control Area Size
Section 2.5.3.1: Selecting the Optimal Percentage of Free Space
Section 2.5.3.2: Altering the Free Space Specification When Loading a Data Set

Rule DAS604: Excessive secondary extents were allocated

Finding: CPExpert determined that an excessive number of secondary extents were allocated for the keyed sequenced data set (KSDS) or variable relative record data set (VRRDS) VSAM data sets listed.

Impact: This finding is used to assess problems or potential problems with VSAM primary and secondary allocation values for the data sets listed. The impact can be significant, particularly with on-line applications.

Discussion: When a VSAM data set is allocated, space allocation amounts normally are specified for both a primary allocation and a secondary allocation. When the primary amount on the first volume is used up, a secondary amount is allocated on that volume by the end-of-volume (EOV) routine, using the amount specified for the secondary allocation.

This space allocation process can be repeated until the volume is out of space or until the extent limit is reached. Depending on the type of data set allocation request, a new volume may be used if the current volume is out of space.

A large number of I/O operations is involved when the secondary allocation takes place. Consequently, a small amount of space should not be specified for the primary or secondary allocation value, especially for a KSDS data set or for a VRRDS data set¹.

CPExpert examines the SMF Type 64 information contained in MXG TYPE64 data set to identify VSAM KSDS or VRRDS data sets that have excessive secondary allocations.

CPExpert compares NREXTNTS variable (the number of secondary extents in the VSAM data set this OPEN) with the **EXTENTS** guidance variable in USOURCE(DASGUIDE). CPExpert produces Rule DAS604 when the NREXTENT (the total number of extents) is greater than one, and the number of secondary extents allocated for this OPEN exceeds the value specified by the EXTENTS guidance variable.

The default value for the EXTENTS guidance variable is 0, indicating that CPExpert should produce Rule DAS604 when any secondary extent was allocated for the VSAM data sets listed.

¹ If the VSAM data set is used **primarily by on-line applications**, you might establish a relatively small primary allocation (which would result in a small Control Area), so the impact of Control Area splits would be spread across the transactions encountering the Control Area splits. See Rule DAS600 for further discussion about this issue.

The following example illustrates the output from Rule DAS604:

Suggestion: If CPExpert produces Rule DAS604, you should consider the following alternatives:

RULE DAS604: EXCESSIVE SECONDARY EXTENTS WERE ALLOCATED

VOLSER: RLS014. More than 0 secondary extents were allocated for the VSAM data sets listed below. There are a large number of I/O operations involved when the secondary allocation takes place. Depending on your I/O configuration and constraints, a large number of secondary allocations could cause significant performance degradation of applications. A large number of secondary extents normally means that the secondary allocation value is too small. This problem can be particularly serious for KSDS or VRRDS data sets.

SMF TIME STAMP	JOB NAME	VSAM DATA SET	TOTAL EXTENTS	EXTENTS THIS OPEN	FILE TYPE
10:00,29AUG2000	CICS2AGC	RLSADSW.VF03D.DATAENDB.DATA.....	10	6	KSDS DATA
10:30,29AUG2000	CICS2AGC	RLSADSW.VF03D.DATAENDB.DATA.....	10	6	KSDS DATA

- If a relatively large number of secondary allocations is made, you should consider increasing the primary allocation amount for the data set.
- Additionally, you should consider increasing the secondary allocation amount for the data set.
- If the above actions are not appropriate, you can change the EXTENTS guidance variable in USOURCE(DASGUIDE). Section 3 describes how to change the EXTENTS guidance variable if you feel that Rule DAS604 is produced too often, or if you do not wish to take action when secondary extents are allocated.
- Alternatively, you can exclude the reported VSAM data sets from analysis. Section 3 describes how to exclude VSAM data sets from analysis. However, you should be aware that no analysis of potential VSAM problems will be performed on data sets that are excluded from analysis.

Reference: *DFSMS: Using Data Sets* (SC26-7339 for OS/390; SC26-7410 for z/OS)
Section 2.2.2.4: Allocating Space for VSAM Data Sets

VSAM Demystified Redbook (SG24-6105)
Section 2.6: Parameters affecting performance

Rule DAS605: Excessive extents were used and secondary allocation was small

Finding: CPExpert determined that an excessive number of secondary extents were allocated for the keyed sequenced data set (KSDS) or a variable relative record data set (VRRDS) VSAM data sets listed.

Impact: This finding is used to assess problems or potential problems with VSAM primary and secondary allocation values for the data sets listed. The impact can be significant, particularly with on-line applications.

Discussion: When a VSAM data set is allocated, space allocation amounts normally are specified for both a primary allocation and a secondary allocation. When the primary amount on the first volume is used up, a secondary amount is allocated on that volume by the end-of-volume (EOV) routine, using the amount specified for the secondary allocation.

This space allocation process can be repeated until the volume is out of space or until the extent limit is reached. Depending on the type of data set allocation request, a new volume may be used if the current volume is out of space.

A large number of I/O operations is involved when the secondary allocation takes place. Consequently, a small amount of space should not be specified for the primary or secondary allocation value, especially for a KSDS data set or for a VRRDS data set¹.

CPExpert examines the SMF Type 64 information contained in MXG TYPE64 data set to identify VSAM KSDS or VRRDS data sets that have excessive secondary allocations.

CPExpert compares NREXTNTS variable (the number of secondary extents in the VSAM data set this OPEN) with the **EXTENTS** guidance variable in USOURCE(DASGUIDE). CPExpert produces Rule DAS605 when the NREXTENT (the total number of extents) is greater than one, and the number of secondary extents allocated for the current OPEN exceeds the value specified by the EXTENTS guidance variable.

When Rule DAS604 is produced, CPExpert analyzes the primary and secondary allocation units. If the primary or secondary allocation unit is in tracks, CPExpert produces Rule DAS605 to reflect the allocation values.

¹ If the VSAM data set is used **primarily by on-line applications**, you might establish a relatively small primary allocation (which would result in a small Control Area), so the impact of Control Area splits would be spread across the transactions encountering the Control Area splits. See Rule DAS600 for further discussion about this issue.

The following example illustrates the output from Rule DAS605:

RULE DAS605: PRIMARY OR SECONDARY ALLOCATION UNIT WAS SMALL

VOLSER: RLS01C. More than 0 secondary extents were allocated for the VSAM data sets listed below. Either the primary allocation or the secondary allocation (or both) was smaller than one cylinder. In addition to causing multiple extents, this allocation size also means that the Control Area (CA) size is less than one cylinder. You should consider increasing the allocation units to use cylinders rather than tracks. The below shows the number of tracks used for the space allocation, and the number of extents:

SMF TIME STAMP	JOB NAME	VSAM DATA SET	..	PRIMARY ALLOCATION	-----SECONDARY----- ALLOCATION	EXTENTS
10:30,29AUG2000	CICS2ACA	RLSADSW.VF04D.DATAENDB.INDEX.....	..	2 TRKS	1 TRK	3
11:00,29AUG2000	CICS2ACA	RLSADSW.VF04D.DATAENDB.INDEX.....	..	2 TRKS	1 TRK	3

Suggestion: If CPExpert produces Rule DAS605, you should consider the following alternatives:

- If a relatively large number of secondary allocations is made, you should consider increasing the primary allocation amount for the data set.
- Alternatively, you should consider increasing the secondary allocation amount for the data set.
- If the above actions are not appropriate, you can change the EXTENTS guidance variable in USOURCE(DASGUIDE). Section 3 describes how to change the EXTENTS guidance variable if you feel that Rule DAS600 is produced too often, or if you do not wish to take action when secondary extents are allocated.
- Alternatively, you can exclude the reported VSAM data sets from analysis. Section 3 describes how to exclude VSAM data sets from analysis. However, you should be aware that no analysis of potential VSAM problems will be performed on data sets that are excluded from analysis.

Reference: *DFSMS: Using Data Sets* (SC26-7339 for OS/390; SC26-7410 for z/OS)
Section 2.2.2.4: Allocating Space for VSAM Data Sets

VSAM Demystified Redbook (SG24-6105)
Section 2.6: Parameters affecting performance

Rule DAS606: Primary or Secondary allocation unit was small

Finding: CPExpert determined that an excessive number of secondary extents were allocated for the keyed sequenced data set (KSDS) or variable relative record data set (VRRDS) VSAM data sets listed, and either the primary or secondary allocation unit was small.

Impact: This finding is used to assess problems or potential problems with VSAM primary and secondary allocation values for the data sets listed. The impact can be significant, particularly with on-line applications.

Discussion: When a VSAM data set is allocated, space allocation amounts normally are specified for both a primary allocation and a secondary allocation. When the primary amount on the first volume is used up, a secondary amount is allocated on that volume by the end-of-volume (EOV) routine, using the amount specified for the secondary allocation.

This space allocation process can be repeated until the volume is out of space or until the extent limit is reached. Depending on the type of data set allocation request, a new volume may be used if the current volume is out of space.

A large number of I/O operations is involved when the secondary allocation takes place. Consequently, a small amount of space should not be specified for the primary or secondary allocation value, especially for a KSDS data set or for a VRRDS data set¹.

CPExpert examines the SMF Type 64 information contained in MXG TYPE64 data set to identify VSAM KSDS or VRRDS data sets that have excessive secondary allocations.

CPExpert compares NREXTNTS variable (the number of secondary extents in the VSAM data set this OPEN) with the **EXTENTS** guidance variable in USOURCE(DASGUIDE). CPExpert produces Rule DAS604 when the NREXTENT (the total number of extents) is greater than one, and the number of secondary extents allocated for this OPEN exceeds the value specified by the EXTENTS guidance variable.

¹ If the VSAM data set is used **primarily by on-line applications**, you might establish a relatively small primary allocation (which would result in a small Control Area), so the impact of Control Area splits would be spread across the transactions encountering the Control Area splits. See Rule DAS600 for further discussion about this issue.

When Rule DAS604 is produced, CPExpert analyzes the primary and secondary allocation units. If the primary or secondary allocation unit is in cylinders, CPExpert produces Rule DAS606 to reflect the allocation values.

The following example illustrates the output from Rule DAS606:

RULE DAS606: PRIMARY OR SECONDARY ALLOCATION UNIT WAS SMALL

VOLSER: RLS014. More than 0 secondary extents were allocated for the VSAM data sets listed below, even though cylinders were used as the allocation unit. The allocation specified for either the primary or secondary allocation should be increased. The below shows the primary and secondary space allocation and number of extents:

SMF TIME STAMP	JOB NAME	VSAM DATA SET	PRIMARY ALLOCATION	SECONDARY ALLOCATION	EXTENTS
10:00,29AUG2000	CICS2AGC	RLSADSW.VF03D.DATAENDB.DATA.....	30 CYLS	10 CYLS	9
10:30,29AUG2000	CICS2AGC	RLSADSW.VF03D.DATAENDB.DATA.....	30 CYLS	10 CYLS	9

Suggestion: If CPExpert produces Rule DAS606, you should consider the following alternatives:

- If a relatively large number of secondary allocations were made, you should consider increasing the primary allocation amount for the data set.
- Alternatively, you should consider increasing the secondary allocation amount for the data set.
- If the above actions are not appropriate, you can change the EXTENTS guidance variable in USOURCE(DASGUIDE). Section 3 describes how to change the EXTENTS guidance variable if you feel that Rule DAS600 is produced too often, or if you do not wish to take action when excessive secondary extents are allocated.
- Alternatively, you can exclude the reported VSAM data sets from analysis. Section 3 describes how to exclude VSAM data sets from analysis. However, you should be aware that no analysis of potential VSAM problems will be performed on data sets that are excluded from analysis.

Reference: *DFSMS: Using Data Sets* (SC26-7339 for OS/390; SC26-7410 for z/OS)
Section 2.2.2.4: Allocating Space for VSAM Data Sets

VSAM Demystified Redbook (SG24-6105)
Section 2.6: Parameters affecting performance

Rule DAS607: VSAM data set is close to maximum number of extents

Finding: CPExpert determined that a VSAM data set has had a significant number of secondary allocations, and is close to reaching the maximum number of extents.

Impact: This finding is used to assess potential problems with VSAM data sets reaching the maximum number of extents. The finding can have a HIGH IMPACT if the VSAM data set reaches the maximum number of extents.

Discussion: When a VSAM data set is allocated, space allocation amounts normally are specified for both a primary allocation and a secondary allocation. When the primary amount on the first volume is used up, a secondary amount is allocated on that volume by the end-of-volume (EOV) routine, using the amount specified for the secondary allocation.

This space allocation process can be repeated until the volume is out of space or until the extent limit is reached. Depending on the type of data set allocation request, a new volume may be used if the current volume is out of space.

A VSAM data set can have up to 255 extents per component, and a striped VSAM data set can have up to 255 extents per stripe. The last four extents are reserved for extending a data set when the last extent cannot be allocated in one piece. (VSAM attempts to extend a data set when the total number of extents is less than 250.

When a multivolume VSAM data set extends to the next volume, the data class specifies if the initial space allocated on that volume is the primary or secondary amount. The default is the primary amount. After the primary amount of space is used up, space is allocated in secondary amounts. By using a data class, it is possible to indicate whether to take a primary or secondary amount when VSAM extends to a new volume.

An ABEND results when an extend is attempted, but the maximum number of extents was reached. CPExpert attempts to provide an “early warning” of this potential situation.

CPExpert examines the SMF Type 64 information contained in MXG TYPE64 data set to identify VSAM data sets that have used a significant number of extents, such that there is danger of reaching the maximum extents. CPExpert compares NREXTENT variable (the total number of extents in the VSAM data set) with the **MXEXTENT** guidance variable in USOURCE(DASGUIDE). CPExpert produces Rule DAS607 when the

NREXTENT is greater than the value specified by the MXEXTENT guidance variable **and** at least one extent was allocated during the current OPEN of the data set (CPEXpert uses the NREXTENTS variable in TYPE64 for this decision).

The default value of the MXEXTENT guidance variable is 225, indicating that CPEXpert should produce Rule DAS607 when at least 225 extents have been allocated for a VSAM data set. Since the maximum allowable is 255, the default value provides a threshold at which CPEXpert provides notification that there is a potential problem.

The following example illustrates the output from Rule DAS607:

RULE DAS607: VSAM DATA SET IS CLOSE TO MAXIMUM NUMBER OF EXTENTS

VOLSER: RLS003. More than 225 extents were allocated for the VSAM data sets listed below. The VSAM data sets are approaching the maximum number of extents allowed. The below shows the number of extents and the primary and secondary space allocation:

SMF TIME STAMP	JOB NAME	VSAM DATA SET	TOTAL EXTENTS	EXTENTS THIS OPEN	---ALLOCATIONS---
					PRIMARY SECONDARY
10:30,11MAR2002	CICS2ABA	RLSADSW.VF01D.DATAENDB.DATA.....	229	4	30 CYL 100 CYL

Suggestion: If CPEXpert produces Rule DAS607, you should consider the following alternatives:

- Determine (potentially by consulting with applications personnel) whether the VSAM data set is expected to continue increasing. If the VSAM data set is expected to continue increasing, you should take action. IBM suggests that you (1) use the access method services REPRO command to make a backup copy of the cluster that contains the data set, (2) delete the cluster from the catalog with the DELETE command, (3) use the DEFINE command to redefine the cluster in the catalog with increased space allocation, and (4) reload the backup of the cluster with the REPRO command. If this alternative is selected, you should consider **significantly** increasing the primary allocation value and consider increasing the secondary allocation value.
- Alternatively, you should examine the amount of space actually used in the VSAM data set by logical records, compared with free space or space used by deleted logical records (depending on the type of VSAM data set). It is possible that inappropriate values have been specified for the FREESPACE parameter.

The percentages of free space should yield full records and full control intervals with a minimum amount of unusable space. If too much free space was specified for the Control Intervals or Control Areas, there

would be more direct access storage required to contain the data set. This unnecessary free space could cause the data set to be unnecessarily large, and cause extents to be acquired more often.

You can use LISTCAT to examine statistics about the space used and the free space allocated and used.

- If the above actions are not appropriate, you can change the MXEXTENT guidance variable in USOURCE(DASGUIDE). Section 3 describes how to change the MXEXTENT guidance variable if you feel that Rule DAS607 is produced prematurely.
- Alternatively, you can exclude the reported VSAM data sets from analysis. Section 3 describes how to exclude VSAM data sets from analysis. However, you should be aware that no analysis of potential VSAM problems will be performed on data sets that are excluded from analysis.

Reference: *DFSMS: Using Data Sets* (SC26-7339 for OS/390; SC26-7410 for z/OS)
Section 2.2.2.4: Allocating Space for VSAM Data Sets
Section 2.5.3: Optimizing Free Space Distribution

VSAM Demystified Redbook (SG24-6105)
Section 2.6: Parameters affecting performance

Rule DAS610: Relatively small CI size was used for sequential processing

Finding: CPExpert noticed that a relatively small CI size was used for sequential processing. This finding applies only if SMF Type 42 (Data Set Statistics)¹ and SMF Type 64 (VSAM Statistics) records are available in a MXG performance data base.

Impact: This can have a LOW IMPACT, MEDIUM IMPACT, or HIGH IMPACT on the performance of applications referencing the VSAM data. The level of impact depends on the number of sequential I/O operations.

Discussion: A VSAM file structure consists of one or more *Control Intervals (CIs)* and one or more *Control Areas (CAs)*.

- A **Control Interval** is a continuous area of direct access storage that VSAM uses to store logical records. When a logical record is read from direct access storage, the entire Control Interval containing the record is read into a VSAM buffer in virtual storage. The desired logical record is then transferred from the VSAM buffer to a user-defined buffer or work area. While logical records within a Control Interval may vary in length, all Control Intervals in a specific VSAM data set are of the same length.
- A **Control Area** contains one or more Control Intervals. The Control Intervals are grouped together into fixed-length contiguous areas of direct access storage. A VSAM data set is composed of one or more Control Areas.

Each VSAM data set is defined as a cluster of one or more components.

- The *data component* is the part of a VSAM data set, alternate index, or catalog that contains the data records. The minimum size of a Control Area for a data component is one track, and the maximum size is one cylinder of DASD storage.
- The *index component* is a collection of logically sequenced keys. A key is a value taken from a fixed defined field in each logical record in the VSAM data set. The key identifies the record's position in the data set. Using the index, VSAM is able to randomly retrieve a record from the data component when a request is made for a record with a certain key. The size of the Control Area for an index component is one track of DASD storage.

¹%LET TYPE42DS = Y; must be specified in USOURCE(GENGUIDE) must be specified in USOURCE(ENGLIST) or USOURCE(DASGUIDE) to advise CPExpert that TYPE42DS is available.

The size of Control Intervals can vary from one VSAM data set to another, but all the Control Intervals of a particular data set component must be the same length. Control interval size affects record processing speed and storage requirements:

- Data sets with large control interval sizes require more buffer space in virtual storage.
- Data sets with large control interval sizes require fewer I/O operations to bring a given number of records into virtual storage, because fewer index records must be read.
- Free space is used more efficiently as control interval size increases relative to data record size. This is because there are fewer Control Interval splits and less wasted space.

The type of processing that is used should guide the choice of control interval size, particularly for the data component:

- **Sequential processing.** When sequential processing accounts for most of the accesses, a large Control Interval for the data component would normally be a good choice. This is because multiple records can be read into buffers and processed sequentially. For example, given a 16KB data buffer space, it is better to read two 8 KB Control Intervals with one I/O operation, than to read four 4 KB Control Intervals with two I/O operations.
- **Direct processing.** When direct processing accounts for most of the accesses, a small Control Interval for the data component is preferable. This is because only one logical record is retrieved at a time with direct access. Since a Control Interval normally contains several logical records, a large Control Interval would create unnecessary I/O overhead reading the logical records that would not be accessed. IBM suggests that a 4096 byte Control Interval normally would be appropriate for direct access.
- **Mixed processing.** If the processing is a mixture of sequential and direct, a small Control Interval for the data component with multiple buffers for sequential processing can be a good choice. The small Control Interval would reduce the I/O processing when using direct access to reference a logical record. The multiple buffers would allow read-ahead (and I/O efficiency) when using sequential processing.

After applying the screening criteria specified for VSAM data sets, and extracting SMF Type 64 information for those VSAM data sets, CPExpert examines SMF Type 42 (Data Set Statistics) information for the selected VSAM data sets. CPExpert uses the TYPE42DS information to compute

the percent of sequential accesses to the data component of the VSAM data set, using the following algorithm:

$$\text{Percent sequential accesses} = \frac{S42AMSRB}{S42AMSRB + S42AMDRB}$$

where: S42AMSRB = Blocks read using sequential access

S42AMDRB = Blocks read using direct access

CPEXpert produces Rule DAS610 when the percent of sequential accesses for the data component was greater than the **PCTSEQ** guidance variable in USOURCE(DASGUIDE), and the Control Interval size was less than 8192 bytes.

The default value for the PCTSEQ guidance variable is 80%, so CPEXpert will produce Rule DAS610 when more than 80% of the accesses were sequential for the data component, and the Control Interval size was less than 8192 bytes.

The following example illustrates the output from Rule DAS610:

RULE DAS610: RELATIVELY SMALL CI SIZE USED FOR SEQUENTIAL PROCESSING

VOLSER: D83NE2. A relatively small data Control Interval (CI) size was used for the VSAM data sets listed below, yet more than 80% of the accesses were for sequential processing. You normally should define a CI size of 8KB or larger for VSAM data sets used mostly for sequential processing. The I/O RATE is for the time the data set was open.

SMF TIME STAMP	JOB NAME	VSAM DATA SET	I/O RATE	-ACCESS TYPE (PCT)-		CI SIZE
				SEQUENTIAL	DIRECT	
17:00,14OCT1997	NETAK31.	OPSYS.NV31.J90.AOC.DSILOGP.DATA.....	19.2	100.0	0.0	4096

While not shown in the above example, CPEXpert also shows the OPEN time for the VSAM data set. This normally is the duration of the current OPEN. If %LET VSAMSMRY=Y; was specified in USOURCE(DASGUIDE), the OPEN time represents the sum of the times the VSAM data was OPEN for all TYPE64 records in the performance data base.

Suggestion: If CPEXpert produces Rule DAS610, you should consider the following alternatives:

- For sequential processing of logical records, system performance normally will be better if multiple logical records are read with a single I/O operation. Consequently, a larger Control Interval would normally improve performance since read-ahead I/O processing could be performed, and you should consider increasing the Control Interval.

-
- If the above action is not appropriate, you can change the **PCTSEQ** guidance variable in USOURCE(DASGUIDE). Section 3 describes how to change the PCTSEQ guidance variable if you feel that Rule DAS610 is produced prematurely.
 - Alternatively, you can exclude the reported VSAM data sets from analysis. Section 3 describes how to exclude VSAM data sets from analysis. However, you should be aware that no analysis of potential VSAM problems will be performed on data sets that are excluded from analysis.

Reference: *DFSMS: Using Data Sets* (SC26-7339 for OS/390; SC26-7410 for z/OS)
Section 2.5.1.2: Data Control Interval Size

Rule DAS611: Relatively large CI size was used for direct processing

Finding: CPExpert noticed that a relatively large CI size was used for direct processing. This finding applies only if SMF Type 42 (Data Set Statistics)¹ and SMF Type 64 (VSAM Statistics) records are available in a MXG performance data base.

Impact: This can have a LOW IMPACT, MEDIUM IMPACT, or HIGH IMPACT on the performance of applications referencing the VSAM data. The level of impact depends on the number of direct I/O operations.

Discussion: A VSAM file structure consists of one or more *Control Intervals (CIs)* and one or more *Control Areas (CAs)*.

- A **Control Interval** is a continuous area of direct access storage that VSAM uses to store logical records. When a logical record is read from direct access storage, the entire Control Interval containing the record is read into a VSAM buffer in virtual storage. The desired logical record is then transferred from the VSAM buffer to a user-defined buffer or work area. While logical records within a Control Interval may vary in length, all Control Intervals in a specific VSAM data set are of the same length.
- A **Control Area** contains one or more Control Intervals. The Control Intervals are grouped together into fixed-length contiguous areas of direct access storage. A VSAM data set is composed of one or more control areas.

Each VSAM data set is defined as a cluster of one or more components.

- The *data component* is the part of a VSAM data set, alternate index, or catalog that contains the data records. The minimum size of a Control Area for a data component is one track, and the maximum size is one cylinder of DASD storage.
- The *index component* is a collection of logically sequenced keys. A key is a value taken from a fixed defined field in each logical record in the VSAM data set. The key identifies the record's position in the data set. Using the index, VSAM is able to randomly retrieve a record from the data component when a request is made for a record with a certain key. The size of the Control Area for an index component is one track of DASD storage.

¹%LET TYPE42DS = Y; must be specified in USOURCE(GENGUIDE) must be specified in USOURCE(GENGUIDE) or USOURCE(DASGUIDE) to advise CPExpert that TYPE42DS is available.

The size of Control Intervals can vary from one VSAM data set to another, but all the Control Intervals of a particular data set component must be the same length. Control interval size affects record processing speed and storage requirements:

- Data sets with large control interval sizes require more buffer space in virtual storage.
- Data sets with large control interval sizes require fewer I/O operations to bring a given number of records into virtual storage, because fewer index records must be read.
- Free space is used more efficiently as control interval size increases relative to data record size. This is because there are fewer Control Interval splits and less wasted space.

The type of processing that is used should guide the choice of control interval size, particularly for the data component:

- **Sequential processing.** When sequential processing accounts for most of the accesses, a large Control Interval for the data component would normally be a good choice. This is because multiple records can be read into buffers and processed sequentially. For example, given a 16KB data buffer space, it is better to read two 8 KB Control Intervals with one I/O operation, than to read four 4 KB Control Intervals with two I/O operations.
- **Direct processing.** When direct processing accounts for most of the accesses, a small Control Interval for the data component is preferable. This is because only one logical record is retrieved at a time with direct access. Since a Control Interval normally contains several logical records, a large Control Interval would create unnecessary I/O overhead reading the logical records that would not be accessed. IBM suggests that a 4096 byte Control Interval normally would be appropriate for direct access.
- **Mixed processing.** If the processing is a mixture of sequential and direct, a small Control Interval for the data component with multiple buffers for sequential processing can be a good choice. The small Control Interval would reduce the I/O processing when using direct access to reference a logical record. The multiple buffers would allow read-ahead (and I/O efficiency) when using sequential processing.

After applying the screening criteria specified for VSAM data sets, and extracting SMF Type 64 information for those VSAM data sets, CPExpert examines SMF Type 42 (Data Set Statistics) information for the selected VSAM data sets. CPExpert uses the TYPE42DS information to compute

the percent of accesses to the data component of the VSAM data set that were direct, using the following algorithm:

$$\text{Percent direct accesses} = \frac{S42AMDRB}{S42AMSRB + S42AMDRB}$$

where: S42AMSRB = Blocks read using sequential access

S42AMDRB = Blocks read using direct access

CPEXpert produces Rule DAS611 when the percent of direct accesses for the data component was greater than the **PCTDIR** guidance variable in USOURCE(DASGUIDE), and the Control Interval size was greater than 4096 bytes. Additionally, CPEXpert verifies that the maximum logical records (maximum LRECL) is less than 50% of the Control Interval size. This verification is done to make sure that Rule DAS611 is not produced for VSAM data sets that have spanned records.

The default value for the PCTDIR guidance variable is 80%, so CPEXpert will produce Rule DAS611 when more than 80% of the accesses were direct for the data component, and the Control Interval size was more than 4096 bytes.

The following example illustrates the output from Rule DAS611:

RULE DAS611: RELATIVELY LARGE CI SIZE USED FOR DIRECT PROCESSING

VOLSER: CICS0E. A relatively large data Control Interval (CI) size was used for the VSAM data sets listed below, yet more than 80% of the accesses were for direct processing. You normally should define a CI size of 4KB for VSAM data sets used mostly for direct processing. The I/O RATE is for the time the data set was open.

SMF TIME STAMP	JOB NAME	VSAM DATA SET	..	I/O	-ACCESS TYPE (PCT)-	CI	
			..	RATE	SEQUENTIAL	DIRECT	SIZE
17:00,14OCT1997	DFHCOMDS	CICS410.PET.DFHCS.D.DATA.....		174.7	0.0	100.0	18432

While not shown in the above example, CPEXpert also shows the OPEN time for the VSAM data set. This normally is the duration of the current OPEN. If %LET VSAMSMRY=Y; was specified in USOURCE(DASGUIDE), the OPEN time represents the sum of the times the VSAM data was OPEN for all TYPE64 records in the performance data base.

Suggestion: If CPEXpert produces Rule DAS611, you should consider the following alternatives:

-
- Large Control Interval sizes for direct access normally degrade performance because the entire Control Interval (consisting of multiple logical records) is brought into storage with each access to a record. With direct processing, normally there is a low probability that more than one record in the Control Interval will be referenced. Consequently, you should consider decreasing the Control Interval size.
 - If the above action is not appropriate, you can change the **PCTDIR** guidance variable in USOURCE(DASGUIDE). Section 3 describes how to change the PCTSEQ guidance variable if you feel that Rule DAS611 is produced prematurely.
 - Alternatively, you can exclude the reported VSAM data sets from analysis. Section 3 describes how to exclude VSAM data sets from analysis. However, you should be aware that no analysis of potential VSAM problems will be performed on data sets that are excluded from analysis.

Reference: *DFSMS: Using Data Sets* (SC26-7339 for OS/390; SC26-7410 for z/OS)
Section 2.5.1.2: Data Control Interval Size

Rule DAS612: Relatively large CI size was used for mixed processing

Finding: CPExpert noticed that a relatively large CI size was used for processing a mixture of sequential and direct accesses. This finding applies only if SMF Type 42 (Data Set Statistics)¹ and SMF Type 64 (VSAM Statistics) records are available in a MXG performance data base.

Impact: This can have a LOW IMPACT, MEDIUM IMPACT, or HIGH IMPACT on the performance of applications referencing the VSAM data. The level of impact depends on the number of I/O operations.

Discussion: A VSAM file structure consists of one or more *Control Intervals (CIs)* and one or more *Control Areas (CAs)*.

- A **Control Interval** is a continuous area of direct access storage that VSAM uses to store logical records. When a logical record is read from direct access storage, the entire Control Interval containing the record is read into a VSAM buffer in virtual storage. The desired logical record is then transferred from the VSAM buffer to a user-defined buffer or work area. While logical records within a Control Interval may vary in length, all Control Intervals in a specific VSAM data set are of the same length.
- A **Control Area** contains one or more Control Intervals. The Control Intervals are grouped together into fixed-length contiguous areas of direct access storage. A VSAM data set is composed of one or more control areas.

Each VSAM data set is defined as a cluster of one or more components.

- The *data component* is the part of a VSAM data set, alternate index, or catalog that contains the data records. The minimum size of a Control Area for a data component is one track, and the maximum size is one cylinder of DASD storage.
- The *index component* is a collection of logically sequenced keys. A key is a value taken from a fixed defined field in each logical record in the VSAM data set. The key identifies the record's position in the data set. Using the index, VSAM is able to randomly retrieve a record from the data component when a request is made for a record with a certain key. The size of the Control Area for an index component is one track of DASD storage.

¹%LET TYPE42DS = Y; must be specified in USOURCE(GENGUIDE) must be specified in USOURCE(GENGUIDE) or USOURCE(DASGUIDE) to advise CPExpert that TYPE42DS is available.

The size of Control Intervals can vary from one VSAM data set to another, but all the Control Intervals of a particular data set component must be the same length. Control interval size affects record processing speed and storage requirements:

- Data sets with large control interval sizes require more buffer space in virtual storage.
- Data sets with large control interval sizes require fewer I/O operations to bring a given number of records into virtual storage, because fewer index records must be read.
- Free space is used more efficiently as control interval size increases relative to data record size. This is because there are fewer Control Interval splits and less wasted space.

The type of processing that is used should guide the choice of control interval size, particularly for the data component:

- **Sequential processing.** When sequential processing accounts for most of the accesses, a large Control Interval for the data component would normally be a good choice. This is because multiple records can be read into buffers and processed sequentially. For example, given a 16KB data buffer space, it is better to read two 8 KB Control Intervals with one I/O operation, than to read four 4 KB Control Intervals with two I/O operations.
- **Direct processing.** When direct processing accounts for most of the accesses, a small Control Interval for the data component is preferable. This is because only one logical record is retrieved at a time with direct access. Since a Control Interval normally contains several logical records, a large Control Interval would create unnecessary I/O overhead reading the logical records that would not be accessed. IBM suggests that a 4096 byte Control Interval normally would be appropriate for direct access.
- **Mixed processing.** If the processing is a mixture of sequential and direct, a small Control Interval for the data component with multiple buffers for sequential processing can be a good choice. The small Control Interval would reduce the I/O processing when using direct access to reference a logical record. The multiple buffers would allow read-ahead (and I/O efficiency) when using sequential processing.

After applying the screening criteria specified for VSAM data sets, and extracting SMF Type 64 information for those VSAM data sets, CPExpert

examines SMF Type 42 (Data Set Statistics) information for the selected VSAM data sets. CPExpert uses the TYPE42DS information to compute the percent of sequential accesses to the data component of the VSAM data set, using the following algorithm:

$$\text{Percent sequential accesses} = \frac{S42AMSRB}{S42AMSRB + S42AMDRB}$$

where: S42AMSRB = Blocks read using sequential access

S42AMDRB = Blocks read using direct access

CPExpert also uses the TYPE42DS information to compute the percent of accesses to the VSAM data set that were direct, using the following algorithm:

$$\text{Percent direct accesses} = \frac{S42AMDRB}{S42AMS4B + S42AMDRB}$$

CPExpert produces Rule DAS612 when:

- The percent of direct accesses for the data component was greater than 30%, **and**
- the percent of sequential accesses for the data component was greater than 30%, **and**
- the maximum logical record length was less than 50% of the Control Interval size, **and**
- the Control Interval size was greater than 4096 bytes, **and**
- less than 10 buffers had been allocated.

The result of this algorithm will select those VSAM data sets with data component accessed for both sequential and direct accesses, that can have more than one logical record per Control Interval, and that have not been allocated a relatively large number of buffers.

The following example illustrates the output from Rule DAS612:

RULE DAS612: RELATIVELY LARGE CI SIZE USED FOR MIXED (RANDOM AND SEQ)

VOLSER: HOLDE2. A relatively large data Control Interval (CI) size was used for the VSAM data sets listed below and few buffers were defined for these data sets, yet the data sets were used for a mixture of random and sequential processing. You normally should define a CI size of 4KB for these VSAM data sets (for random processing) but define a relatively large number of buffers (for sequential processing). The I/O RATE is for the time the data set was open.

SMF TIME STAMP	JOB NAME	VSAM DATA SET	..	I/O RATE	-ACCESS TYPE (PCT)- SEQUENTIAL	DIRECT	CI SIZE	BUFFERS
17:00,14OCT1997	ARI170022	IMSDSW4P.DSW.RECON1.DATA.....	..	32.4	50.0	50.0	32768	2

While not shown in the above example, CPExpert also shows the OPEN time for the VSAM data set. This normally is the duration of the current OPEN. If %LET VSAMSMRY=Y; was specified in USOURCE(DASGUIDE), the OPEN time represents the sum of the times the VSAM data was OPEN for all TYPE64 records in the performance data base.

Suggestion: If CPExpert produces Rule DAS612, you should consider the following alternatives:

- Decrease the Control Interval size and increase the number of buffers. As mentioned earlier, IBM suggests that a small data Control Interval with multiple buffers for sequential processing can be a good choice with mixed processing.
- If the above action is not appropriate, you can change the PCTSEQ guidance variable in USOURCE(DASGUIDE). Section 3 describes how to change the PCTSEQ guidance variable if you feel that Rule DAS612 is produced prematurely.
- Alternatively, you can exclude the reported VSAM data sets from analysis. Section 3 describes how to exclude VSAM data sets from analysis. However, you should be aware that no analysis of potential VSAM problems will be performed on data sets that are excluded from analysis.

Reference: *DFSMS: Using Data Sets* (SC26-7339 for OS/390; SC26-7410 for z/OS)
Section 2.5.1.2: Data Control Interval Size

Rule DAS620: The number of data buffers should be increased

Finding: CPExpert noticed that Non-Shared resources (NSR) was specified for VSAM data sets and most of the access was sequential processing. However, relatively few data buffers were assigned to the data sets. Consequently, I/O processing was inefficient. This finding applies only if SMF Type 42 (Data Set Statistics)¹ and SMF Type 64 (VSAM Statistics) records are available in a MXG performance data base.

Impact: This can have a LOW IMPACT, MEDIUM IMPACT, or HIGH IMPACT on the performance of applications referencing the VSAM data. The level of impact depends on the number of direct I/O operations that are used.

Discussion: A VSAM file structure consists of one or more *Control Intervals (CIs)* and one or more *Control Areas (CAs)*.

- A **Control Interval** is a continuous area of direct access storage that VSAM uses to store logical records. When a logical record is read from direct access storage, the entire Control Interval containing the record is read into a VSAM buffer in virtual storage. The desired logical record is then transferred from the VSAM buffer to a user-defined buffer or work area. While logical records within a Control Interval may vary in length, all Control Intervals in a specific VSAM data set are of the same length.
- A **Control Area** contains one or more Control Intervals. The Control Intervals are grouped together into fixed-length contiguous areas of direct access storage. A VSAM data set is composed of one or more Control Areas.

I/O buffers are used by VSAM to read and write control intervals from DASD to virtual storage. For a key-sequenced data set (KSDS) or variable-length relative records data sets (VRRDS), VSAM requires a minimum of three buffers: two buffers for data control intervals² and one buffer for an index control interval. Only data buffers are needed for entry-sequenced, for linear data sets (LSDS), and for fixed-length relative record data sets (RRDS).

The VSAM defaults provide these **minimum** buffers. However, to increase performance, there are parameters to override the VSAM default values.

¹%LET TYPE42DS = Y; must be specified in USOURCE(GENGUIDE) must be specified in USOURCE(GENGUIDE) or USOURCE(DASGUIDE) to advise CPExpert that TYPE42DS is available.

²One of the data buffers is used only for formatting control areas and splitting control intervals and control areas.

Selecting good buffering options can reduce the number of I/O operations, reduce job elapsed time, reduce CPU time, device reduce connect time, and reduce device disconnect time. IBM benchmarks³ have shown that, depending on the workload, selecting optimal buffering parameters can provide 90% reduction in the number of I/O operations, provide over 65% reduction in job elapsed, provide over 40% reduction in CPU time, and provide over 50% reduction in device connect time. These remarkable savings were achieved simply by altering the buffering characteristics of jobs processing VSAM data sets.

The optimum buffering techniques vary, depending the buffering techniques (or resource pools) used and how the records are accessed (sequential or direct). There are four types of resource pools, depending on the type of data sharing that is implemented:

- **Non-Shared Resource (NSR).** NSR is the default VSAM buffering technique. With NSR, VSAM buffers are not shared among VSAM data sets, and the buffers are located in the private area. VSAM data sets with NSR buffering can be accessed sequentially or direct (or both). However, NSR is suited for sequential processing because, if the data set access is sequential, the buffers are managed with a read-ahead algorithm. The read-ahead algorithm provides overlap of I/O and CPU processing and is efficient for sequential accesses. Since NSR is oriented toward sequential access, there is no expectation that a record will be re-used (as might exist with direct processing). Consequently, once a record is processed from the NSR buffers, the buffer is likely to be reclaimed for another record read from DASD.
- **Local Shared Resource (LSR).** With LSR, VSAM buffers normally are shared among VSAM data sets accessed by tasks in the same address space. Since LSR is oriented toward shared (and direct) access, there is an expectation that a record might be re-used. Consequently, buffer management algorithms retain buffers as long as possible, using a least-recently used (LRU) algorithm, after a record is processed from the LSR buffers. There is no read-ahead algorithm with LSR, and there is no inherent overlap of I/O and CPU processing. LSR is appropriate for direct access of VSAM data sets, regardless of whether the data sets are shared.
- **Global Shared Resource (GSR).** GSR provides serialization of shared resources across multiple systems. In a GSR complex, programs can serialize access to data sets on shared DASD volumes at the *data set level* rather than at the *DASD volume level*. A program on one system can access one data set on a shared volume while other programs on any system can access other data sets on the same volume.

³VSAM Demystified, SG24-6105, Section 2.6.9 (Buffering options)

-
- **Record Level Sharing (RLS).** VSAM RLS provides multisystem sharing of VSAM data sets in a parallel sysplex, using cross-system locking. CICS is the exploiter of RLS. RLS enables VSAM data to be shared, with full update capability, between many applications running in many CICS regions. With RLS, CICS regions that share VSAM data sets can reside in one or more MVS images within an MVS parallel sysplex.

As described above, NSR is best used for applications that use sequential or skip sequential as their primary access mode. When data sets are accessed sequentially, performance can be increased by specifying multiple buffers for the data component⁴. When there are multiple data buffers, VSAM uses a read-ahead function to read the next data control intervals into buffers as buffers become available. On the other hand, more buffers than necessary might cause excessive paging or excessive internal processing. There is an optimum point at which more buffers do not continue to improve performance.

Please note that additional buffers do not improve performance if the VSAM data set is processed with direct access. This is because there is no read-ahead I/O activity.

After applying the screening criteria specified for VSAM data sets, and extracting SMF Type 64 information for those VSAM data sets, CPExpert examines SMF Type 42 (Data Set Statistics) information for the selected VSAM data sets. CPExpert uses the TYPE42DS information to compute the percent of sequential accesses to the VSAM data set, using the following algorithm:

$$\text{Percent sequential accesses} = \frac{S42AMSRB}{S42AMSRB + S42AMDRB}$$

where: S42AMSRB = Blocks read using sequential access
S42AMDRB = Blocks read using direct access

CPExpert produces Rule DAS620 when the TYPE42DS S42DSBUF variable showed that NSR was used, the percent of sequential accesses for the data component was greater than 75%, and less than 10 buffers had been assigned to the data component for the VSAM data set.

⁴ Having only one I/O buffer for the index component does not hinder performance, because VSAM gets to the next CI by using the horizontal pointers in sequence set records rather than the vertical pointers in the index set. Extra index buffers have little effect during sequential processing.

The following example illustrates the output from Rule DAS620:

RULE DAS620: THE NUMBER OF DATA BUFFERS SHOULD BE INCREASED

VOLSER: D83NE2. Non-Shared resources (NSR) was specified for the below VSAM data sets and most of the access was sequential processing. However, relatively few data buffers were assigned to the data sets. You should consider increasing the number of data buffers to at least 10 buffers, and preferably up to 30 buffers. The I/O RATE is for the time the data set was open.

SMF TIME STAMP	JOB NAME	VSAM DATA SET	..	I/O RATE	-ACCESS TYPE (PCT)-	..	BUFFERS
					SEQUENTIAL	DIRECT	ASSIGNED
17:00,14OCT1997	NETAK31.	OPSYS.NV31.J90.AOC.DSILOGP.DATA.....		40.2	100.0	0.0	2

While not shown in the above example, CPEXpert also shows the OPEN time for the VSAM data set. This normally is the duration of the current OPEN. If %LET VSAMSMRY=Y; was specified in USOURCE(DASGUIDE), the OPEN time represents the sum of the times the VSAM data was OPEN for all TYPE64 records in the performance data base.

Suggestion: If CPEXpert produces Rule DAS620, you should consider the following alternatives:

- Increase the number of buffers for the VSAM data sets identified by Rule DAS620. IBM benchmarks⁵ have shown that with read sequential access, performance improves significantly if the number of buffers for the data component of VSAM data sets is increased from the default 2 buffers to 10 buffers. Performance continues to increase as more buffers are added, but the improvement was less in IBM's benchmarks, as the number of buffers increased beyond 30 buffers. Consequently, IBM recommends that the number of buffers for the data component of sequentially accessed VSAM data sets be increased from the default to 10 buffers or up to 30 buffers if space permits. This recommendation is only a general guideline, however. The optimum number of data component buffers varies according to the amount of CPU processing⁶ done between each read to the data set.

⁵VSAM Demystified, SG24-6105, Section 2.6.9 (Buffering options)

⁶The amount of CPU processing can be done by the application accessing the VSAM data set, or can be done by work (either system or application work) executing at an equal or higher CPU dispatching priority. The CPU-I/O overlap efficiencies created by multiple VSAM buffers for NSR sequential access occur only when the application otherwise would wait on logical records. Application delay waiting for I/O can be analyzed by the WLM Component of CPEXpert, or can be seen in RMF reports for the service class to which the application belongs.

There is one significant exception to this “increase buffers” recommendation; that exception occurs when SHAREOPTIONS 4 has been specified for the VSAM data set.

The SHAREOPTIONS parameter specifies how the component or cluster can be shared among users within one system or across systems. With SHAREOPTIONS 4, the data set can be fully shared by any number of users. This setting does not allow any type of non-RLS access when the data set is already open for RLS processing. With this option, each user is responsible for maintaining both read and write integrity for the data. With SHAREOPTIONS 4, buffers are refreshed at each request. Also, the read-ahead function has no effect and defer write is not used. **Therefore, for SHAREOPTIONS 4, keeping data buffers at a minimum can actually improve performance.**

Unfortunately, SMF information does not describe the SHAREOPTIONS value, so CPExpert is unable to detect that SHAREOPTIONS 4 has been specified for a VSAM data set, and thus cannot suppress this finding for VSAM data sets with SHAREOPTIONS 4.

- Alternatively, you can specify System Managed Buffering (SMB)⁷ for the VSAM data sets listed. VSAM can use system-managed buffering to determine the number of buffers and the type of buffer management to use for VSAM data sets. VSAM also determines the number of buffers to locate in Hiperspace for use in direct optimization. To indicate that VSAM is to use SMB, specify either of the following options:
 - Specify the ACCBIAS subparameter of the JCL DD statement AMP parameter and specify **Sequential Optimized (SO)** for the record access bias. This technique provides the most efficient buffers for sequential application. Approximately 500K of processor virtual storage for buffers is required for this technique, defaulted to above 16 MB.

If 500K processor virtual storage for buffers is not available with the application, consider specifying **Sequential Weighted (SW)** for the record access bias. This technique provides efficient buffers for sequential application. Approximately 100K of processor virtual storage for buffers is required for this technique, defaulted to above 16 MB,.

- Specify Record Access Bias in the data class and an application processing option in the ACB.

⁷ Please review *DFSMS: Using Data Sets* (Section 2.5.4.2.3: Processing Guidelines and Restrictions) before implementing system managed buffering.

For system-managed buffering (SMB), the data set must use both of the following options:

- System Management Subsystem (SMS) storage
- Extended format (DSNTYPE=ext in the data class)
- Alternatively, you can exclude the reported VSAM data sets from analysis. Section 3 describes how to exclude VSAM data sets from analysis. However, you should be aware that no analysis of potential VSAM problems will be performed on data sets that are excluded from analysis.

Reference: *DFSMS: Using Data Sets* (SC26-7339 for OS/390; SC26-7410 for z/OS)
Section 2.5.5.5.1 Data Buffers for Sequential Access Using Nonshared Resources
Section 2.5.4.2: Tuning for System-Managed Buffering
Section 2.5.4.2.3: Processing Guidelines and Restrictions

IBM Redbook: *VSAM Demystified* (SG24-6105)
Section 2.6.9 (Buffering options)

Rule DAS621: The number of index buffers should be increased (direct access)

Finding: CPExpert noticed that direct processing accounted for a significant amount of the I/O activity to the index component of VSAM data sets. However, insufficient buffers were assigned to the index component. The number of index buffers should be increased for optimal performance of the VSAM data sets listed. This finding applies only if SMF Type 42 (Data Set Statistics)¹ and SMF Type 64 (VSAM Statistics) records are available in a MXG performance data base.

Impact: This can have a LOW IMPACT, MEDIUM IMPACT, or HIGH IMPACT on the performance of applications referencing the VSAM data. The level of impact depends on the number of direct I/O operations.

Discussion: A VSAM file structure consists of one or more *Control Intervals (CIs)* and one or more *Control Areas (CAs)*.

- A **Control Interval** is a continuous area of direct access storage that VSAM uses to store logical records. When a logical record is read from direct access storage, the entire Control Interval containing the record is read into a VSAM buffer in virtual storage. The desired logical record is then transferred from the VSAM buffer to a user-defined buffer or work area. While logical records within a Control Interval may vary in length, all Control Intervals in a specific VSAM data set are of the same length.
- A **Control Area** contains one or more Control Intervals. The Control Intervals are grouped together into fixed-length contiguous areas of direct access storage. A VSAM data set is composed of one or more Control Areas.

Each VSAM data set is defined as a cluster of one or more components.

- The *data component* is the part of a VSAM data set, alternate index, or catalog that contains the data records. The minimum size of a Control Area for a data component is one track, and the maximum size is one cylinder of DASD storage.
- The *index component* is a collection of logically sequenced keys. A key is a value taken from a fixed defined field in each logical record in the VSAM data set. The key identifies the record's position in the data set. Using the index, VSAM is able to randomly retrieve a record from the

¹%LET TYPE42DS = Y; must be specified in USOURCE(GENGUIDE) must be specified in USOURCE(GENGUIDE) or USOURCE(DASGUIDE) to advise CPExpert that TYPE42DS is available.

data component when a request is made for a record with a certain key. The size of the Control Area for an index component is one track of DASD storage.

Key-sequenced data sets (KSDS) and variable-length relative records data sets (VRRDS) contain both a data component and an index component. Additionally, each alternate index contains a data component and an index component. Entry-sequenced data sets (ESDS), linear data sets (LDS), and fixed-length relative record data sets (RRDS) contain only a data component.

The index component consists of two parts: *sequence set* and *index set*.

- The sequence set is the lowest level of index control intervals and directly points to the data Control Interval in the data Control Area. There is one Control Interval in the sequence set for each *data* Control Area. This *index* Control Interval contains pointers and high key information for each *data* Control Interval. The index Control Interval also contains horizontal pointers from one sequence set Control Interval to the next higher keyed sequence set Control Interval.
- The records in all levels of the index above the sequence set are called the index set. If there is more than one sequence set Control Interval, VSAM automatically builds another index level. An entry in an index set record consists of the highest possible key in an index record in the next lower level, and a pointer to the beginning of that index record.

I/O buffers are used by VSAM to read and write control intervals from DASD to virtual storage. A minimum of one buffer is required for an index Control Interval. Having only one index I/O buffer does not hinder performance when the VSAM data set is accessed sequentially, because VSAM gets to the next Control Interval by using the horizontal pointers in sequence set records rather than the vertical pointers in the index set.

The minimum of one index buffer² is inadequate with **direct** access to the VSAM data set. When using direct access to retrieve a record from a key-sequenced data set or variable-length RRDS (or store a record using keyed access), VSAM needs to examine the index of the data set. When an index record must be retrieved to locate a data record, VSAM makes room for the new index record by deleting the index record that VSAM judges to be least useful under the prevailing circumstances. If only one index buffer were provided, a serious performance problem would occur if an index record were continually deleted from virtual storage to make

²Multiple buffers for the *data component* do not increase performance with direct processing, because only one data buffer is used for each access.

room for another index record, and then retrieved again later when it is required.

Providing more than the minimum number of index can significantly improve performance. IBM benchmarks³ show over 50% reduction in EXCPs and almost 50% reduction in CPU time, by increasing the number of index buffers from 1 to 4 for non-shared resource (NSR) with direct access.

Unused index buffers do not normally degrade performance, so an adequate number should be specified. For optimum performance, the number of index buffers should be at least as large as the number of high-level index set Control Intervals, plus one per string to contain the entire high-level index set and one sequence set control interval per string in virtual storage.

After applying the screening criteria specified for VSAM data sets, and extracting SMF Type 64 information for those VSAM data sets, CPExpert examines SMF Type 42 (Data Set Statistics) information for the selected VSAM data sets. CPExpert uses the TYPE42DS information (MXG variable DSTYPE) to identify KSDS and VRRDS VSAM data sets. CPExpert selects these VSAM data sets that were used with NSR (using MXG variable S42DSBUF). CPExpert uses the TYPE42DS information to compute the percent of direct accesses to the VSAM data set, using the following algorithm:

$$\text{Percent direct accesses} = \frac{S42AMDRB}{S42AMSRB + S42AMDRB}$$

where: S42AMSRB = Blocks read using sequential access

S42AMDRB = Blocks read using direct access

CPExpert produces Rule DAS621 under the following conditions:

- The TYPE42DS S42DSBUF variable showed that NSR was used for KSDS or VRRDS VSAM data sets, **and**
- The percent of direct accesses for the index component was greater than **DIRINDEX** guidance variable in USOURCE(DASGUIDE), **and**
- The number of buffers (the MXG BUFDRNO variable) assigned to the index component was less than the number of index levels (the MXG ACCLEVEL variable) for the VSAM data set.

³VSAM Demystified, SG24-6105, Table 6 (NSR buffering with direct access - STRNO=1)

The default value for the DIRINDEX guidance variable is 25%, so CPExpert will produce Rule DAS621 for NSR VSAM data sets when more than 25% of the accesses were direct for the index component and the number of buffers assigned to the index component was less than the number of index levels.

The following example illustrates the output from Rule DAS621:

RULE DAS621: THE NUMBER OF INDEX BUFFERS SHOULD BE INCREASED

VOLSER: PRD005. Non-Shared resources (NSR) was specified as the buffering technique for the below VSAM data sets, and most of the accesses were random processing. However, relatively few index buffers were assigned to the data sets. You should consider increasing the number of index buffers to the number of index levels. The I/O RATE is for the time the data set was open.

SMF TIME STAMP	JOB NAME	VSAM DATA SET	I/O RATE	-ACCESS TYPE (PCT)- SEQUENTIAL	DIRECT	INDEX LEVELS	BUFFERS ASSIGNED
10:04,19SEP2002	PGRED01D	PWFWFA.PK.GENTRAN.EDIOEA.INDEX.....	30.7	0.0	100.0	3	1
10:20,19SEP2002	PFED01D	PWFWFA.PK.GENTRAN.EDIOEA.INDEX.....	15.7	0.0	100.0	3	1

While not shown in the above example, CPExpert also shows the OPEN time for the VSAM data set. This normally is the duration of the current OPEN. If %LET VSAMSMRY=Y; was specified in USOURCE(DASGUIDE), the OPEN time represents the sum of the times the VSAM data was OPEN for all TYPE64 records in the performance data base.

Suggestion: If CPExpert produces Rule DAS621, you should consider the following alternatives:

- **Increase the number of index buffers.** You should consider increasing the number of index buffers allocated for the VSAM data sets listed by Rule DAS621.
- **Consider the use of System Managed Buffering.** If the data sets listed by Rule DAS621 are SMS managed, have extended format, and your installation has DFSMS V1R4 or later, you can use system managed buffering (SMB)⁴. System managed buffering enables VSAM to determine the optimum number of index and data buffers, as well as the type of buffer management (LSR or NSR). IBM benchmarks⁵ show up to

⁴Please review *DFSMS: Using Data Sets* (Section 2.5.4.2.3: Processing Guidelines and Restrictions) before implementing system managed buffering.

⁵*VSAM Demystified*, SG24-6105, Table 7 (Direct access: Benefits of using SMF: Updates and insertions)

90% reduction in EXCPs, DASD connect time, and CPU time when converting to system managed buffering!

To indicate that VSAM is to use SMB, specify either of the following options:

- Specify the ACCBIAS subparameter of the JCL DD statement AMP parameter.
- Specify Record Access Bias in the data class and an application processing option in the ACB. The ACB must be NSR and MACRF cannot contain any of the following:
 - ICI (Improved control interval processing)
 - AIX Processing the data set through the alternate index of the path specified in the DDname
 - UBF (Management of I/O buffers left up to the VSAM user)

Regardless of whether the JCL DD statement AMP is used or the data class/ACB is used, you should specify the following for the record access bias:

- **Specify Direct Optimized (DO)** for the record access bias for the record access bias if *direct access is near 100% for the VSAM data sets listed*. The DO processing technique optimizes for totally random record access. This technique overrides the user specification for non-shared resources (NSR) buffering with a local shared resources (LSR) implementation of buffering.
- **Specify Direct Weighted (DW)** for the record access bias if the data sets are processed using a mixture of direct and sequential access, and *direct access is dominate*. DW processing provides the minimum read-ahead buffers for sequential retrieval and the maximum index buffers for direct requests.
- **Specify Sequential Weighted (SW)** for the record access bias if the data sets are processed using a mixture of direct and sequential access, and *sequential access is dominate*. This technique uses read-ahead buffers for sequential requests and provides additional index buffers for direct requests. Approximately 100K of processor virtual storage for buffers is required for this technique, defaulted to above 16 MB.

-
- **Consider the use of Batch LSR.** For batch applications **and if Rule DAS621 shows that direct processing⁶ is near 100%**, you can use Batch LSR. Batch LSR provides advantages in an application using VSAM NSR buffering techniques to switch to LSR without changing the application source code or link-editing the application again. Only a JCL change is required.
 - If the above actions are not appropriate, you can change the **DIRINDEX** guidance variable in USOURCE(DASGUIDE). Section 3 describes how to change the DIRINDEX guidance variable if you feel that Rule DAS621 is produced prematurely.
 - Alternatively, you can exclude the reported VSAM data sets from analysis. Section 3 describes how to exclude VSAM data sets from analysis. However, you should be aware that no analysis of potential VSAM problems will be performed on data sets that are excluded from analysis.

Reference: *DFSMS: Using Data Sets* (SC26-7339 for OS/390; SC26-7410 for z/OS)
Section 2.5.5.5.1 Data Buffers for Sequential Access Using Nonshared Resources
Section 2.5.4.2: Tuning for System-Managed Buffering

IBM Redbook: *VSAM Demystified* (SG24-6105)
Section 2.6.9 (Buffering options)

⁶Using the Batch LSR subsystem with sequential access could degrade performance rather than improve it. This is because Batch LSR does not provide "read-ahead" capability. The "read-ahead" capability is essential for good performance with sequential access.

Rule DAS622: The number of index buffers should be increased (STRNO value)

Finding: The number of index buffers was not greater than the number of strings specified in the Access Control Block (ACB) STRNO value. The number of index buffers should be increased for optimal performance of the VSAM data sets listed. This finding applies only if SMF Type 42 (Data Set Statistics)¹ and SMF Type 64 (VSAM Statistics) records are available in a MXG performance data base.

Impact: This can have a LOW IMPACT, MEDIUM IMPACT, or HIGH IMPACT on the performance of applications referencing the VSAM data. The level of impact depends on the number of direct I/O operations.

Discussion: A VSAM file structure consists of one or more *Control Intervals (CIs)* and one or more *Control Areas (CAs)*.

- A **Control Interval** is a continuous area of direct access storage that VSAM uses to store logical records. When a logical record is read from direct access storage, the entire Control Interval containing the record is read into a VSAM buffer in virtual storage. The desired logical record is then transferred from the VSAM buffer to a user-defined buffer or work area. While logical records within a Control Interval may vary in length, all Control Intervals in a specific VSAM data set are of the same length.
- A **Control Area** contains one or more Control Intervals. The Control Intervals are grouped together into fixed-length contiguous areas of direct access storage. A VSAM data set is composed of one or more Control Areas.

Each VSAM data set is defined as a cluster of one or more components.

- The *data component* is the part of a VSAM data set, alternate index, or catalog that contains the data records. The minimum size of a Control Area for a data component is one track, and the maximum size is one cylinder of DASD storage.
- The *index component* is a collection of logically sequenced keys. A key is a value taken from a fixed defined field in each logical record in the VSAM data set. The key identifies the record's position in the data set. Using the index, VSAM is able to randomly retrieve a record from the data component when a request is made for a record with a certain key.

¹%LET TYPE42DS = Y; must be specified in USOURCE(GENGUIDE) must be specified in USOURCE(GENGUIDE) or USOURCE(DASGUIDE) to advise CPEXpert that TYPE42DS is available.

The size of the Control Area for an index component is one track of DASD storage.

Key-sequenced data sets (KSDS) and variable-length relative records data sets (VRRDS) contain both a data component and an index component. Additionally, each alternate index contains a data component and an index component. Entry-sequenced data sets (ESDS), linear data sets (LDS), and fixed-length relative record data sets (RRDS) contain only a data component.

The index component consists of two parts: *sequence set* and *index set*.

- The sequence set is the lowest level of index control intervals and directly points to the data Control Interval in the data Control Area. There is one Control Interval in the sequence set for each *data* Control Area. This *index* Control Interval contains pointers and high key information for each *data* Control Interval. This index Control Interval also contains horizontal pointers from one sequence set Control Interval to the next higher keyed sequence set Control Interval.
- The records in all levels of the index above the sequence set are called the index set. If there is more than one sequence set Control Interval, VSAM automatically builds another index level. An entry in an index set record consists of the highest possible key in an index record in the next lower level, and a pointer to the beginning of that index record.

I/O buffers are used by VSAM to read and write control intervals from DASD to virtual storage. A minimum of one buffer is required for an index Control Interval. Having only one index I/O buffer does not hinder performance when the VSAM data set is accessed sequentially, because VSAM gets to the next Control Interval by using the horizontal pointers in sequence set records rather than the vertical pointers in the index set.

A string is a request to a VSAM data set requiring data set positioning. If different *concurrent* accesses to the same data set are necessary, multiple strings are used. If multiple strings are used, each string requires exclusive control of an index I/O buffer. Therefore, the value specified for the STRNO parameter (in the ACB or GENCB macro, or AMP parameter) is the minimum number of index I/O buffers required when requests that require concurrent positioning are issued.

If the number of I/O buffers provided for index records is greater than the number of requests that require concurrent positioning (specified by the STRNO parameter), one buffer is used for the highest-level index record. Any additional buffers are used, as required, for other index-set index records. With direct access, the minimum number of index buffers should be one more than the value of the STRNO if VSAM is to keep the

highest-level index record resident. Keeping the highest-level index record resident can significantly improve performance, at a modest increase in virtual storage used for index buffers

For optimum performance, the number of index buffers should be at least as large as the number of high-level index set Control Intervals, **plus one per string to contain the entire high-level index set and one sequence set control interval per string in virtual storage.**

After applying the screening criteria specified for VSAM data sets, and extracting SMF Type 64 information for those VSAM data sets, CPExpert examines SMF Type 42 (Data Set Statistics) information for the selected VSAM data sets. CPExpert uses the TYPE42DS information to compute the percent of accesses to the VSAM data set that were CPExpert uses the TYPE42DS information to compute the percent of direct accesses to the VSAM data set, using the following algorithm:

$$\text{Percent direct accesses} = \frac{S42AMDRB}{S42AMSRB + S42AMDRB}$$

where: S42AMSRB = Blocks read using sequential access

S42AMDRB = Blocks read using direct access

CPExpert produces Rule DAS622 under the following conditions:

- The TYPE42DS S42DSBUF variable showed that NSR was used for KSDS or VRRDS VSAM data sets, **and**
- The percent of direct accesses for the index component was greater than **DIRINDEX** guidance variable in USOURCE(DASGUIDE), **and**
- The STRNO specification in the ACB (the MXG ACBSTRNO) was greater than one (indicating that concurrent accesses had been specified for direct processing), **and**
- The number of buffers (the MXG BUFDRNO variable) assigned to the index component was less than the ACBSTRNO value, plus 1.

The default value for the DIRINDEX guidance variable is 25%, so CPExpert will produce Rule DAS622 for NSR VSAM data sets when more than 25% of the accesses were direct for the index component and the number of buffers assigned to the index component was less than the STRNO value, plus 1.

The following example illustrates the output from Rule DAS622:

RULE DAS622: THE NUMBER OF INDEX BUFFERS SHOULD BE INCREASED (STRNO VALUE)

VOLSER: MVS902. Non-Shared resources (NSR) was specified as the buffering technique for the below VSAM data sets, and most of the access were direct processing. However, relatively few index buffers were assigned to the data sets. You should consider increasing the number of index buffers to 1 more than the number of strings specified in the STRNO in the ACB. The I/O RATE is for the time the data set was open.

SMF TIME STAMP	JOB NAME	VSAM DATA SET	I/O RATE	-ACCESS TYPE (PCT)- SEQUENTIAL	DIRECT	STRINGS	BUFFERS ASSIGNED
10:10,19SEP2002	NDMMON..	SDPDPA.PK.MVSP.RT.NDMTCF.INDEX.....	39.2	0.0	100.0	5	5
10:10,19SEP2002	PCGORDER	SDPDPA.PK.MVSP.RT.NDMDRD.INDEX.....	12.1	0.0	100.0	3	3
10:10,19SEP2002	NDMMON..	SDPDPA.PK.MVSP.RT.NDMTCF.INDEX.....	242.8	0.0	100.0	5	5
10:14,19SEP2002	PGDOS06R	SDPDPA.PK.MVSP.CS.NDMDRD.INDEX.....	12.9	0.0	100.0	3	3
10:14,19SEP2002	PGDOS06R	SDPDPA.PK.MVSP.CS.NDMDRD.INDEX.....	16.5	0.0	100.0	3	3

While not shown in the above example, CPExpert also shows the OPEN time for the VSAM data set. This normally is the duration of the current OPEN. If %LET VSAMSMRY=Y; was specified in USOURCE(DASGUIDE), the OPEN time represents the sum of the times the VSAM data was OPEN for all TYPE64 records in the performance data base.

Suggestion: If CPExpert produces Rule DAS622, you should consider the following alternatives:

- **Increase the number of buffers.** You should consider increasing the number of index buffers for the VSAM data sets identified by Rule DAS622 to be one greater than the number of concurrent accesses (as specified by the STRNO parameter).
- Alternatively, you can specify System Managed Buffering (SMB)² for the VSAM data set. VSAM can use system-managed buffering to determine the number of buffers and the type of buffer management to use for VSAM data sets. VSAM also determines the number of buffers to locate in Hiperspace for use in direct optimization. To indicate that VSAM is to use SMB, specify either of the following options:
 - Specify the ACCBIAS subparameter of the JCL DD statement AMP parameter and specify Sequential Optimized (SO) for the record access bias.
 - Specify Record Access Bias in the data class and an application processing option in the ACB.

²Please review *DFSMS: Using Data Sets* (Section 2.5.4.2.3: Processing Guidelines and Restrictions) before implementing system managed buffering.

For system-managed buffering (SMB), the data set must use both of the following options:

- System Management Subsystem (SMS) storage
- Extended format (DSNTYPE=ext in the data class)
- If the above actions are not appropriate, you can change the **DIRINDEX** guidance variable in USOURCE(DASGUIDE). Section 3 describes how to change the DIRINDEX guidance variable if you feel that Rule DAS622 is produced prematurely.
- Alternatively, you can exclude the reported VSAM data sets from analysis. Section 3 describes how to exclude VSAM data sets from analysis. However, you should be aware that no analysis of potential VSAM problems will be performed on data sets that are excluded from analysis.

Reference: *DFSMS: Using Data Sets* (SC26-7339 for OS/390; SC26-7410 for z/OS)
Section 2.5.4.4.2: Index Buffers for Direct Access
Section 2.5.4.2: Tuning for System-Managed Buffering

IBM Redbook: *VSAM Demystified* (SG24-6105)
Section 2.6.9 (Buffering options)

Rule DAS625: NSR was used, but a large percent of the access was direct

Finding: CPExpert noticed that Non-shared resources (NSR) was used as the access method for the VSAM data sets listed, but most of the accesses were direct. This finding applies only if SMF Type 42 (Data Set Statistics) information is available¹ in a MXG performance data base

Impact: This can have a LOW IMPACT, MEDIUM IMPACT, or HIGH IMPACT on the performance of applications referencing the VSAM data. The level of impact depends on the number of direct I/O operations.

Discussion: I/O buffers are used by VSAM to read and write control intervals from DASD to virtual storage. Two buffering techniques are used with VSAM data sets that are accessed only on a local² system: Non-shared resource (NSR) and Local shared resource (LSR).

- **Non-Shared Resource (NSR).** NSR is the default VSAM buffering technique. With NSR, VSAM buffers are not shared among VSAM data sets, and the buffers are located in the private area. VSAM data sets with NSR buffering can be accessed sequentially or direct (or both). However, NSR is suited for sequential processing because, if the data set access is sequential, the buffers are managed with a read-ahead algorithm. The read-ahead algorithm provides overlap of I/O and CPU processing and is efficient for sequential accesses. Since NSR is oriented toward sequential access, there is no expectation that a record will be re-used (as might exist with direct processing). Consequently, once a record is processed from the NSR buffers, the buffer is likely to be reclaimed for another record read from DASD.
- **Local Shared Resource (LSR).** With LSR, VSAM buffers normally are shared among VSAM data sets accessed by tasks in the same address space. Since LSR is oriented toward shared (and direct) access, there is an expectation that a record might be re-used. Consequently, buffer management algorithms retain buffers as long as possible, using a least-recently used (LRU) algorithm, after a record is processed from the LSR buffers. There is no read-ahead algorithm with LSR, and there is no inherent overlap of I/O and CPU processing. LSR is appropriate for

¹%LET TYPE42DS = Y; must be specified in USOURCE(GENGUIDE) must be specified in USOURCE(GENGUIDE) or USOURCE(DASGUIDE) to advise CPExpert that TYPE42DS is available.

²Two other techniques can be used: global shared resource (GSR) and record level sharing (RLS). GSR provides serialization of shared resources across multiple systems. VSAM RLS provides multisystem sharing of VSAM data sets in a parallel sysplex.

direct access of VSAM data sets, regardless of whether the data sets are shared.

As described above, NSR is best used for applications that use sequential or skip sequential as their primary access mode. NSR is not suited for direct processing, although NSR often **is** used for direct processing because it is easy to use and is the default buffering technique. Nonetheless, performance can be significantly improved if LSR is used for direct processing of VSAM data sets.

After applying the screening criteria specified for VSAM data sets, and extracting SMF Type 64 information for those VSAM data sets, CPExpert examines SMF Type 42 (Data Set Statistics) information for the selected VSAM data sets. CPExpert uses the TYPE42DS information to compute the percent of direct accesses to the VSAM data set, using the following algorithm:

$$\text{Percent direct accesses} = \frac{S42AMDRB}{S42AMSRB + S42AMDRB}$$

where: S42AMSRB = Blocks read using sequential access

S42AMDRB = Blocks read using direct access

CPExpert produces Rule DAS625 under the following conditions:

- The TYPE42DS S42DSBUF variable showed that NSR was used for KSDS or VRRDS VSAM data sets, and
- The percent of direct accesses for the data component was greater than the value specified for the NSRDIR guidance variable in USOURCE(DASGUIDE).

The default value for the NSRDIR guidance variable is 75%, so CPExpert will produce Rule DAS625 when NSR was specified as the buffering technique, and more than 75% of the accesses were direct for the data component.

The following example illustrates the output from Rule DAS625:

RULE DAS625: NSR WAS USED, BUT LARGE PERCENT OF ACCESS WAS DIRECT

VOLSER: MVS902. Non-Shared resources (NSR) was specified as the buffering technique for the below VSAM data sets, but more than 75% of the I/O activity was direct access. NSR is not designed for direct access, and many of the advantages of NSR are not available for direct access. You should consider Local Shared Resources (LSR) for the below VSAM data sets (perhaps using System Managed Buffers to facilitate the use of LSR). The I/O RATE is for the time the data set was open. The SMF TIME STAMP and JOB NAME are from the last record for the data set.

SMF TIME STAMP	JOB NAME	VSAM DATA SET	I/O RATE	OPEN DURATION	-ACCESS TYPE (PCT)- SEQUENTIAL	- DIRECT
13:19,19SEP2002	NRXX807.	SDPDPA.PK.MVSP.RT.NDMGIX.DATA.....	8.4	0:07:08	0.0	100.0
13:19,19SEP2002	NRXX807.	SDPDPA.PR.MVSP.RT.NDMGIXD.DATA.....	11.2	0:06:42	0.0	100.0
13:33,19SEP2002	TSJHM...	SDPDPA.PR.MVSP.RT.NDMRQFDA.DATA.....	0.3	2:21:58	0.0	100.0
13:33,19SEP2002	TSJHM...	SDPDPA.PR.MVSP.RT.NDMRQF.DATA.....	2.8	3:37:53	0.0	100.0
13:33,19SEP2002	TSJHM...	SDPDPA.PK.MVSP.RT.NDMTCF.DATA.....	11.1	6:24:10	0.1	99.9

While not shown in the above example, CPExpert also shows the OPEN time for the VSAM data set. This normally is the duration of the current OPEN. If %LET VSAMSMRY=Y; was specified in USOURCE(DASGUIDE), the OPEN time represents the sum of the times the VSAM data was OPEN for all TYPE64 records in the performance data base.

Suggestion: If CPExpert consistently produces Rule DAS625 for the same VSAM data sets, and the related jobs are having performance problems randomly accessing VSAM data, *you can improve performance with no changes in your applications*. You should consider the following alternatives:

- **Consider the use of System Managed Buffering.** If the data sets listed by Rule DAS625 are SMS managed, have extended format, and your installation has DFSMS V1R4 or later, you can use system managed buffering³. System managed buffering enables VSAM to determine the optimum number of index and data buffers, as well as the type of buffer management (LSR or NSR). IBM benchmarks⁴ show up to 90% reduction in EXCPs, DASD connect time, and CPU time when converting to system managed buffering!

To indicate that VSAM is to use SMB, specify either of the following options:

- Specify the ACCBIAS subparameter of the JCL DD statement AMP parameter.

³Please review *DFSMS: Using Data Sets* (Section 2.5.4.2.3: Processing Guidelines and Restrictions) before implementing system managed buffering.

⁴*VSAM Demystified*, SG24-6105, Table 7 (Direct access: Benefits of using SMF: Updates and insertions)

-
- Specify Record Access Bias in the data class and an application processing option in the ACB. The ACB must be NSR and MACRF cannot contain any of the following:
 - ICI (Improved control interval processing)
 - AIX Processing the data set through the alternate index of the path specified in the DDname
 - UBF (Management of I/O buffers left up to the VSAM user)
 - **Specify Direct Optimized (DO).** Regardless of whether the JCL DD statement AMP is used or the data class/ACB is used, you should specify Direct Optimized (DO) for the record access bias if *direct access is near 100% for the VSAM data sets listed*. The DO processing technique optimizes for totally random record access. This technique overrides the user specification for nonshared resources (NSR) buffering with a local shared resources (LSR) implementation of buffering.
 - **Specify Direct Weighted (DW)** for the record access bias if the data sets are processed using a mixture of direct and sequential access, but direct access is not near⁵ 100%. DW processing provides the minimum read-ahead buffers for sequential retrieval and the maximum index buffers for direct requests.
 - **Consider the use of Batch LSR.** For batch applications **and if Rule DAS625 shows that direct processing⁶ is near 100%**, you can use Batch LSR. Batch LSR provides advantages in an application using VSAM NSR buffering techniques to switch to LSR without changing the application source code or link-editing the application again. Only a JCL change is required.
 - If the above actions are not appropriate, you can change the NSRDIR guidance variable in USOURCE(DASGUIDE). Section 3 describes how to change the NSRDIR guidance variable if you feel that Rule DAS625 is produced prematurely.
 - Alternatively, you can exclude the reported VSAM data sets from analysis. Section 3 describes how to exclude VSAM data sets from analysis. However, you should be aware that no analysis of potential

⁵Note that the default value for the NSRDIR guidance variable is 75, indicating that more than 75% of the processing for the VSAM data set was direct access. Unless this guidance is reduced, Rule DAS625 should not be produced unless more than 75% of the accesses were direct.

⁶Using the Batch LSR subsystem with sequential access could degrade performance rather than improve it. This is because Batch LSR does not provide "read-ahead" capability. The "read-ahead" capability is essential for good performance with sequential access.

VSAM problems will be performed on data sets that are excluded from analysis.

Reference: *DFSMS: Using Data Sets* (SC26-7339 for OS/390; SC26-7410 for z/OS)
Section 2.5.4: Determining I/O Buffer Space for Nonshared Resource
Section 2.5.4.2.3: Processing Guidelines and Restrictions

Rule DAS635: LSR was used, but a large percent of the access was sequential

Finding: CPExpert noticed that local shared resources (LSR) was used as the access method for the VSAM data sets listed, but most of the accesses were sequential. This finding applies only if SMF Type 42 (Data Set Statistics) information is available¹ in a MXG performance data base

Impact: This can have a LOW IMPACT, MEDIUM IMPACT, or HIGH IMPACT on the performance of applications referencing the VSAM data. The level of impact depends on the number of sequential I/O operations.

Discussion: I/O buffers are used by VSAM to read and write control intervals from DASD to virtual storage. Two buffering techniques are used with VSAM data sets that are accessed only on a local² system: Non-shared resource (NSR) and Local shared resource (LSR).

- **Non-Shared Resource (NSR).** NSR is the default VSAM buffering technique. With NSR, VSAM buffers are not shared among VSAM data sets, and the buffers are located in the private area. VSAM data sets with NSR buffering can be accessed sequentially or direct (or both). However, NSR is suited for sequential processing because, if the data set access is sequential, the buffers are managed with a read-ahead algorithm. The read-ahead algorithm provides overlap of I/O and CPU processing and is efficient for sequential accesses. Since NSR is oriented toward sequential access, there is no expectation that a record will be re-used (as might exist with direct processing). Consequently, once a record is processed from the NSR buffers, the buffer is likely to be reclaimed for another record read from DASD.
- **Local Shared Resource (LSR).** With LSR, VSAM buffers normally are shared among VSAM data sets accessed by tasks in the same address space. Since LSR is oriented toward shared (and direct) access, there is an expectation that a record might be re-used. Consequently, buffer management algorithms retain buffers as long as possible, using a least-recently used (LRU) algorithm, after a record is processed from the LSR buffers. There is no read-ahead algorithm with LSR, and there is no inherent overlap of I/O and CPU processing. LSR is appropriate for

¹%LET TYPE42DS = Y; must be specified in USOURCE(GENGUIDE) must be specified in USOURCE(GENGUIDE) or USOURCE(DASGUIDE) to advise CPExpert that TYPE42DS is available.

²Two other techniques can be used: global shared resource (GSR) and record level sharing (RLS). GSR provides serialization of shared resources across multiple systems. VSAM RLS provides multisystem sharing of VSAM data sets in a parallel sysplex.

direct access of VSAM data sets, regardless of whether the data sets are shared.

Selecting good buffering options can reduce the number of I/O operations, reduce job elapsed time, reduce CPU time, device reduce connect time, and reduce device disconnect time. IBM benchmarks³ have shown that, depending on the workload, selecting optimal buffering parameters can provide 90% reduction in the number of I/O operations, provide over 65% reduction in job elapsed, provide over 40% reduction in CPU time, and provide over 50% reduction in device connect time. These remarkable savings were achieved simply by altering the buffering characteristics of jobs processing VSAM data sets.

As described above, **NSR** is best used for applications that use sequential or skip sequential as their primary access mode. When data sets are accessed sequentially, performance can be increased by specifying multiple data buffers. When there are multiple data buffers, VSAM uses a read-ahead function to read the next data control intervals into buffers as buffers become available. On the other hand, more buffers than necessary might cause excessive paging or excessive internal processing. There is an optimum point at which more buffers do not continue to improve performance.

LSR is not suited for applications that use sequential or skip sequential as their primary access mode, because there is no read-ahead algorithm with LSR, and there is no inherent overlap of I/O and CPU processing. Consequently, using LSR for sequential access processing could degrade rather than improve performance.

After applying the screening criteria specified for VSAM data sets, and extracting SMF Type 64 information for those VSAM data sets, CPExpert examines SMF Type 42 (Data Set Statistics) information for the selected VSAM data sets. CPExpert uses the TYPE42DS information to compute the percent of sequential accesses to the VSAM data set, using the following algorithm:

$$\text{Percent sequential accesses} = \frac{S42AMSRB}{S42AMSRB + S42AMDRB}$$

where: S42AMSRB = Blocks read using sequential access

S42AMDRB = Blocks read using direct access

CPExpert produces Rule DAS635 under the following conditions:

³VSAM Demystified, SG24-6105, Section 2.6.9 (Buffering options)

- The TYPE42DS S42DSBUF variable showed that LSR was used for KSDS or VRRDS VSAM data sets, and
- The percent of sequential accesses for the data component was greater than the value specified for the LSRSEQ guidance variable in USOURCE(DASGUIDE).

The default value for the LSRSEQ guidance variable is 75%, so CPExpert will produce Rule DAS635 when LSR was specified as the buffering technique, and more than 75% of the accesses were sequential for the data component.

The following example illustrates the output from Rule DAS635:

RULE DAS635: LSR WAS USED, BUT LARGE PERCENT OF ACCESS WAS SEQUENTIAL

VOLSER: PRD001. Local shared resources (LSR) was specified as the buffering technique for the VSAM data sets listed below, but more than 75% of the I/O activity was sequential access. LSR is not designed for sequential access, and many of the advantages of LSR are not available for sequential access. You should consider Non-Shared Resources (NSR) for the below VSAM data sets. The I/O RATE is for the time the data set was open. The SMF TIME STAMP and JOB NAME are from the last record for the data set.

SMF TIME STAMP	JOB NAME	VSAM DATA SET	I/O RATE	OPEN DURATION	-ACCESS TYPE (PCT)- SEQUENTIAL	DIRECT
10:36,19SEP2002	PICBUYD1	PNPFFA.PK.INFOREM.IC.ROD.DATA.....	1.6	0:02:13	100.0	0.0

While not shown in the above example, CPExpert also shows the OPEN time for the VSAM data set. This normally is the duration of the current OPEN. If %LET VSAMSMRY=Y; was specified in USOURCE(DASGUIDE), the OPEN time represents the sum of the times the VSAM data was OPEN for all TYPE64 records in the performance data base.

Suggestion: If CPExpert consistently produces Rule DAS635 for the same VSAM data sets, and the related jobs are having performance problems accessing VSAM data, you should consider the following alternatives:

- Convert the VSAM buffer technique from LSR to Non-shared resource (NSR) for the VSAM data sets. This will allow VSAM to manage the buffers with a read-ahead algorithm. The read-ahead algorithm will provide overlap of I/O and CPU processing and is efficient for sequential accesses.
- If the above action is not appropriate, you can change the LSRSEQ guidance variable in USOURCE(DASGUIDE). Section 3 describes how to change the LSRSEQ guidance variable if you feel that Rule DAS635 is produced prematurely.

-
- Alternatively, you can exclude the reported VSAM data sets from analysis. Section 3 describes how to exclude VSAM data sets from analysis. However, you should be aware that no analysis of potential VSAM problems will be performed on data sets that are excluded from analysis.

Reference: IBM Redbook: *VSAM Demystified* (SG24-6105)
Section 2.6.9 (Buffering options)

Your turn:

This manual has described how to use the DASD Component to analyze performance constraints with your overall computing environment.

We would appreciate receiving any comments you have regarding this document (style, content, clarity, etc.), or suggestions for improving the DASD Component (ease of use, new rules, changes to rules, etc.). Please send your comments to:

Don Deese
Computer Management Sciences, Inc.
6076-D Franconia Road
Alexandria, VA 22310
www.cpexpert.com

Comments: